# Final Exam Review

# CS 4319

# Data Mining & Warehouses



University of Houston-Downtown

**Fall 2016**

# Data Mining Tasks

- **Prediction Methods**
  - Use some variables to predict unknown or future values of other variables.

  e.g. house prices, medical studies, dep./ind. variables

- **Description Methods**
  - Find human-interpretable patterns that describe the data.

  e.g.: correlation, clusters, trajectories, anomalies)

From [Fayyad, et.al.] Advances in Knowledge Discovery and Data Mining, 1996

# Data Mining Tasks...

- Classification [Predictive]

- Clustering [Descriptive]

- Association Rule Discovery [Descriptive]

- Sequential Pattern Discovery [Descriptive]

- Regression [Predictive]

- Deviation Detection [Predictive]

# Classification: Definition

- Given a collection of records (*training set* )
  - Each record contains a set of *attributes*, one of the attributes is the *class*.
- Find a *model* for class attribute as a function of the values of other attributes.
- Goal: <u>previously unseen</u> records should be assigned a class as accurately as possible.
  - A *test set* is used to determine the accuracy of the model. Usually, the given data set is divided into training and test sets, with training set used to build the model and test set used to validate it.
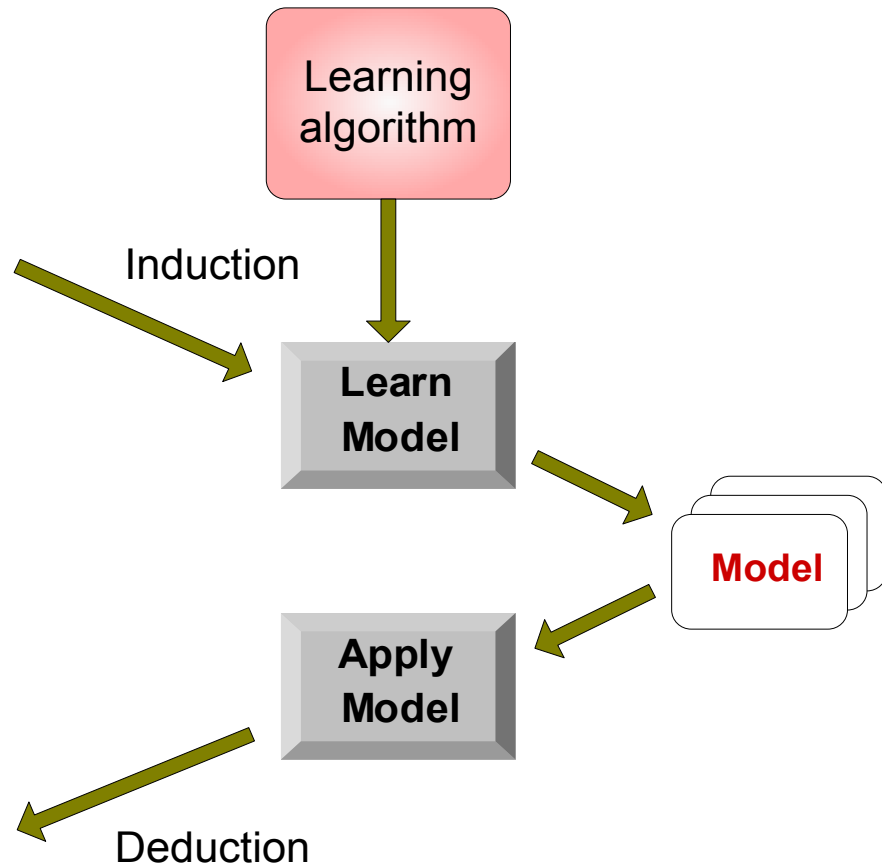
# Illustrating Classification Task

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 1 | Yes | Large | 125K | No |
| 2 | No | Medium | 100K | No |
| 3 | No | Small | 70K | No |
| 4 | Yes | Medium | 120K | No |
| 5 | No | Large | 95K | Yes |
| 6 | No | Medium | 60K | No |
| 7 | Yes | Large | 220K | No |
| 8 | No | Small | 85K | Yes |
| 9 | No | Medium | 75K | No |
| 10 | No | Small | 90K | Yes |

Training Set

Learning algorithm

Induction

Learn Model

Model

| Tid | Attrib1 | Attrib2 | Attrib3 | Class |
|-----|---------|---------|---------|-------|
| 11 | No | Small | 55K | ? |
| 12 | Yes | Medium | 80K | ? |
| 13 | Yes | Large | 110K | ? |
| 14 | No | Small | 95K | ? |
| 15 | No | Large | 67K | ? |

Test Set

Apply Model

Deduction

# Example of a Decision Tree

categorical
categorical
continuous
class

| Tid | Refund | Marital Status | Taxable Income | Cheat |
|-----|--------|----------------|----------------|-------|
| 1 | Yes | Single | 125K | No |
| 2 | No | Married | 100K | No |
| 3 | No | Single | 70K | No |
| 4 | Yes | Married | 120K | No |
| 5 | No | Divorced | 95K | Yes |
| 6 | No | Married | 60K | No |
| 7 | Yes | Divorced | 220K | No |
| 8 | No | Single | 85K | Yes |
| 9 | No | Married | 75K | No |
| 10 | No | Single | 90K | Yes |

**Training Data**

*Splitting Attributes*

Refund
Yes → NO
No → MarSt

MarSt
Single, Divorced → TaxInc
Married → NO

TaxInc
< 80K → NO
> 80K → YES

**Model:  Decision Tree**

# Tree Induction

- Greedy strategy.
    - Split the records based on an attribute test that optimizes certain criterion.

- Issues
    - Determine how to split the records
        - How to specify the attribute test condition?
        - How to determine the best split?
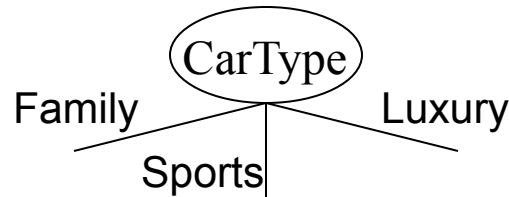    - Determine when to stop splitting

# How to Specify Test Condition?

- Depends on attribute types
  - Nominal
  - Ordinal
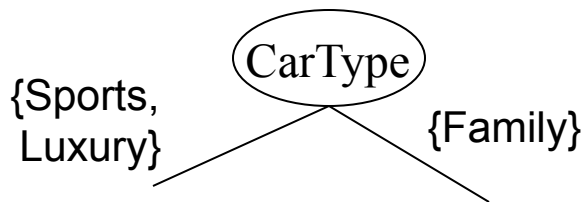  - Continuous

- Depends on number of ways to split
  - 2-way split
  - Multi-way split

# Splitting Based on Nominal Attributes

- **Multi-way split:** Use as many partitions as distinct values.



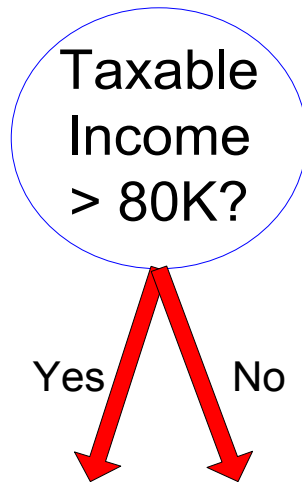- **Binary split:** Divides values into two subsets. Need to find optimal partitioning.
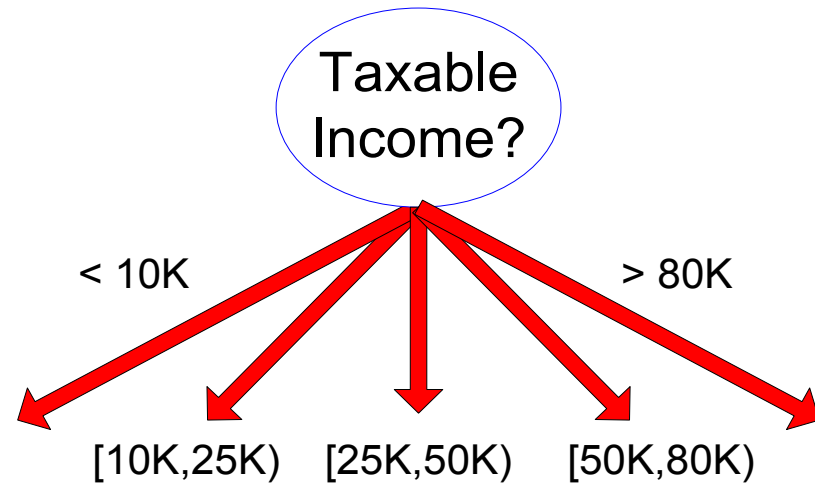
# Splitting Based on Continuous Attributes

- Different ways of handling
  - Discretization to form an ordinal categorical attribute
    - Static – discretize once at the beginning
    - Dynamic – ranges can be found by equal interval bucketing, equal frequency bucketing (percentiles), or clustering.

  - Binary Decision: $(A < v)$ or $(A \geq v)$
    - consider all possible splits and finds the best cut
    - can be more compute intensive

# Splitting Based on Continuous Attributes

Taxable Income > 80K?

Yes    No

(i) Binary split

Taxable Income?

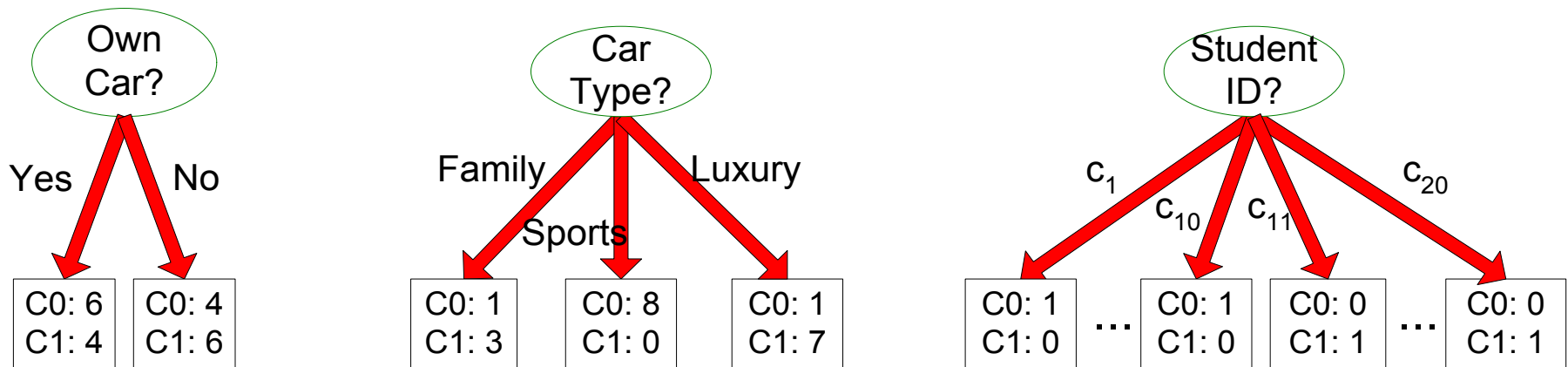< 10K    [10K,25K)    [25K,50K)    [50K,80K)    > 80K

(ii) Multi-way split

# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

# How to determine the Best Split

**Before Splitting: 10 records of class 0,**
**10 records of class 1**



**Which test condition is the best?**

# How to determine the Best Split

● Greedy approach:

   – Nodes with <span style="color:red">homogeneous</span> class distribution are preferred

● Need a measure of node impurity:

| C0: 5 |
|-------|
| C1: 5 |

**Non-homogeneous,**

**High degree of impurity**

| C0: 9 |
|-------|
| C1: 1 |

**Homogeneous,**

**Low degree of impurity**

# Measures of Node Impurity

- Gini Index

- Entropy

- Misclassification error

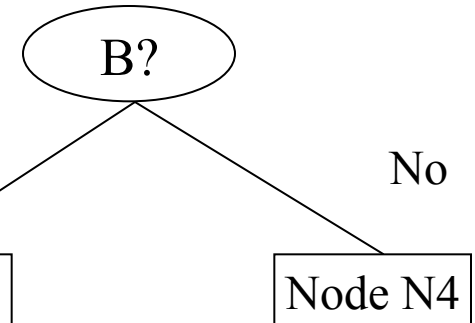# How to Find the Best Split

**Before Splitting:**

| C0 | **N00** |
|----|---------|
| C1 | **N01** |

→ **M0**

**A?**

Yes — No

Node N1 — Node N2

| C0 | **N10** |
|----|---------|
| C1 | **N11** |

| C0 | **N20** |
|----|---------|
| C1 | **N21** |

**M1** — **M2**

**M12**

**B?**

Yes — No

Node N3 — Node N4

| C0 | **N30** |
|----|---------|
| C1 | **N31** |

| C0 | **N40** |
|----|---------|
| C1 | **N41** |

**M3** — **M4**

**M34**

**Gain = M0 – M12 vs  M0 – M34**

# Measure of Impurity: GINI

● Gini Index for a given node t :

$$GINI(t) = 1 - \sum_{j} [p(j \mid t)]^2$$

(NOTE: $p(j \mid t)$ is the relative frequency of class j at node t).

- Maximum ($1 - 1/n_c$) when records are equally distributed among all classes, implying least interesting information

- Minimum (0.0) when all records belong to one class, implying most interesting information

| C1 | 0 |
|----|---|
| C2 | 6 |
| Gini=0.000 | |

| C1 | 1 |
|----|---|
| C2 | 5 |
| Gini=0.278 | |

| C1 | 2 |
|----|---|
| C2 | 4 |
| Gini=0.444 | |

| C1 | 3 |
|----|---|
| C2 | 3 |
| Gini=0.500 | |

# Examples for computing GINI

$$GINI(t) = 1 - \sum_j [p(j \mid t)]^2$$

| C1 | 0 |
|----|---|
| C2 | 6 |

P(C1) = 0/6 = 0    P(C2) = 6/6 = 1

Gini = 1 – P(C1)$^2$ – P(C2)$^2$ = 1 – 0 – 1 = 0

| C1 | 1 |
|----|---|
| C2 | 5 |

P(C1) = 1/6        P(C2) = 5/6

Gini = 1 – (1/6)$^2$ – (5/6)$^2$ = 0.278

| C1 | 2 |
|----|---|
| C2 | 4 |

P(C1) = 2/6        P(C2) = 4/6

Gini = 1 – (2/6)$^2$ – (4/6)$^2$ = 0.444

# Splitting Based on GINI

- Used in CART, SLIQ, SPRINT.
- When a node p is split into k partitions (children), the quality of split is computed as,

$$GINI_{split} = \sum_{i=1}^{k} \frac{n_i}{n} GINI(i)$$

where,      $n_i$ = number of records at child i,

  $n$  = number of records at node p.

# Tree Induction

- Greedy strategy.
  - Split the records based on an attribute test that optimizes certain criterion.

- Issues
  - Determine how to split the records
    - How to specify the attribute test condition?
    - How to determine the best split?
  - Determine when to stop splitting

# Stopping Criteria for Tree Induction

- Stop expanding a node when all the records belong to the same class

- Stop expanding a node when all the records have similar attribute values

- Early termination (to be discussed later)

# Decision Tree Based Classification

- Advantages:
  - Inexpensive to construct
  - Extremely fast at classifying unknown records
  - Easy to interpret for small-sized trees
  - Accuracy is comparable to other classification techniques for many simple data sets

# Example: C4.5

- Simple depth-first construction.

- Uses Information Gain

- Sorts Continuous Attributes at each node.

- Needs entire data to fit in memory.

- Unsuitable for Large Datasets.

  – Needs out-of-core sorting.


- You can download the software from: http://www.cse.unsw.edu.au/~quinlan/c4.5r8.tar.gz

# Practical Issues of Classification

- Underfitting and Overfitting
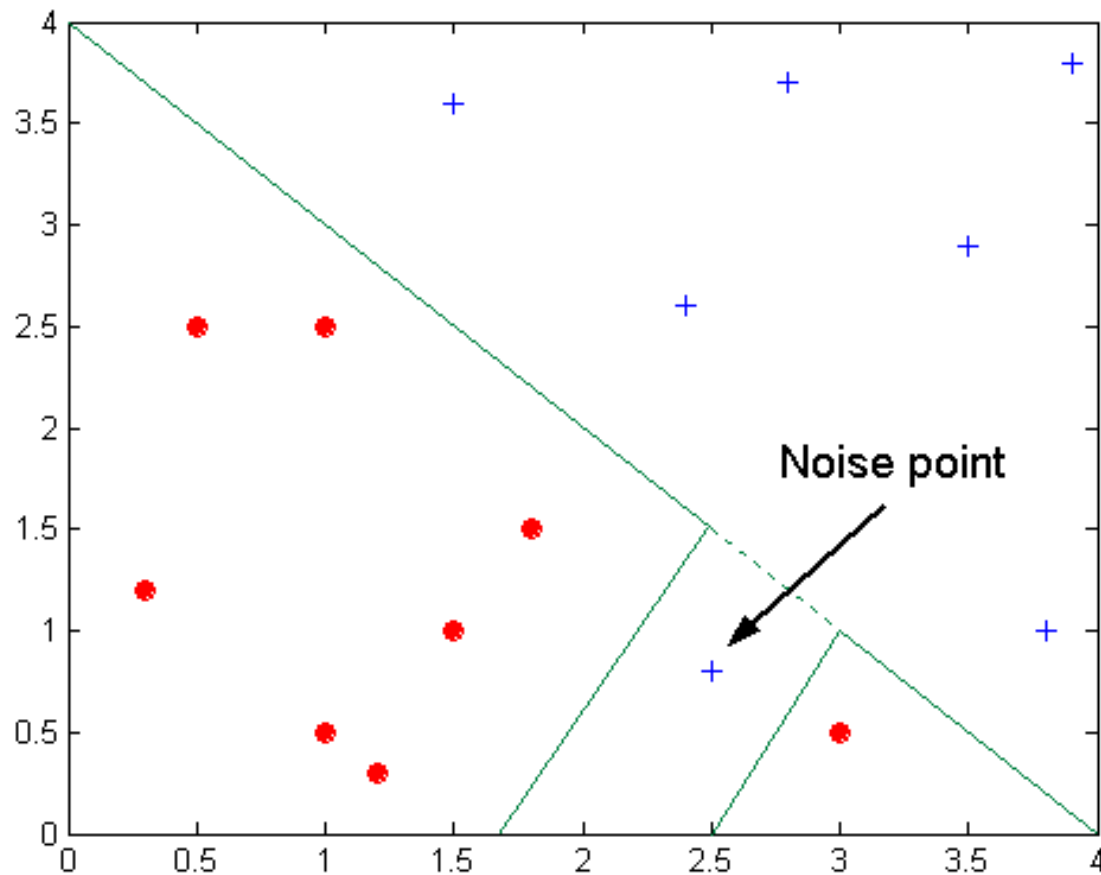
- Missing Values

- Costs of Classification

# Underfitting and Overfitting



**Underfitting**: when model is too simple, both training and test errors are large

# Overfitting due to Noise



**Decision boundary is distorted by noise point**

# Overfitting due to Noise

**Table 4.3.** An example training set for classifying mammals. Class labels with asterisk symbols represent mislabeled records.

| Name | Body Temperature | Gives Birth | Four-legged | Hibernates | Class Label |
|------|------------------|-------------|-------------|------------|-------------|
| porcupine | warm-blooded | yes | yes | yes | yes |
| cat | warm-blooded | yes | yes | no | yes |
| bat | warm-blooded | yes | no | yes | no* |
| whale | warm-blooded | yes | no | no | no* |
| salamander | cold-blooded | no | yes | yes | no |
| komodo dragon | cold-blooded | no | yes | no | no |
| python | cold-blooded | no | no | yes | no |
| salmon | cold-blooded | no | no | no | no |
| eagle | warm-blooded | no | no | no | no |
| guppy | cold-blooded | yes | no | no | no |

* Bats and Whales are misclassified; non-mammals instead of mammals.

# Overfitting due to Noise



(a) Model M1                    (b) Model M2

**Figure 4.25.** Decision tree induced from the data set shown in Table 4.3.

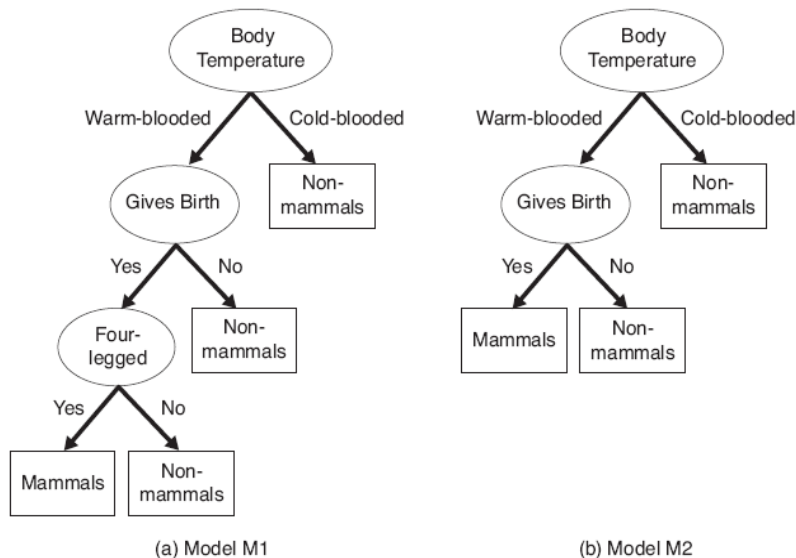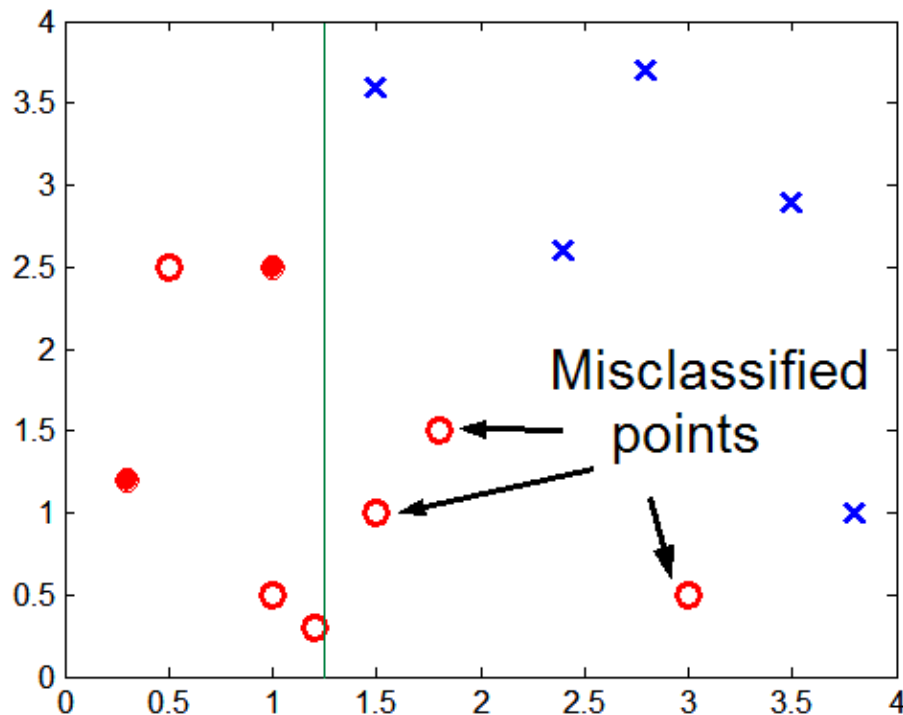# Overfitting due to Noise



(a) Model M1  (b) Model M2

Figure 4.25. Decision tree induced from the data set shown in Table 4.3.

- Decision tree perfectly fits training data (training error=0)

- But error rate on test data is 30%.

- Both humans and dolphins were misclassified as n0n-mammals b/c Body Temp, Gives_Birth and Four-legged values are identical to mislabeled records in training set.

- Spiny anteaters represent an exceptional case (every warm-blooded with no gives_birth is non-mammal in TR_Set

# Overfitting due to Insufficient Examples



**Lack of data points in the lower half of the diagram makes it difficult to predict correctly the class labels of that region**

**- Insufficient number of training records in the region causes the decision tree to predict the test examples using other training records that are irrelevant to the classification task**
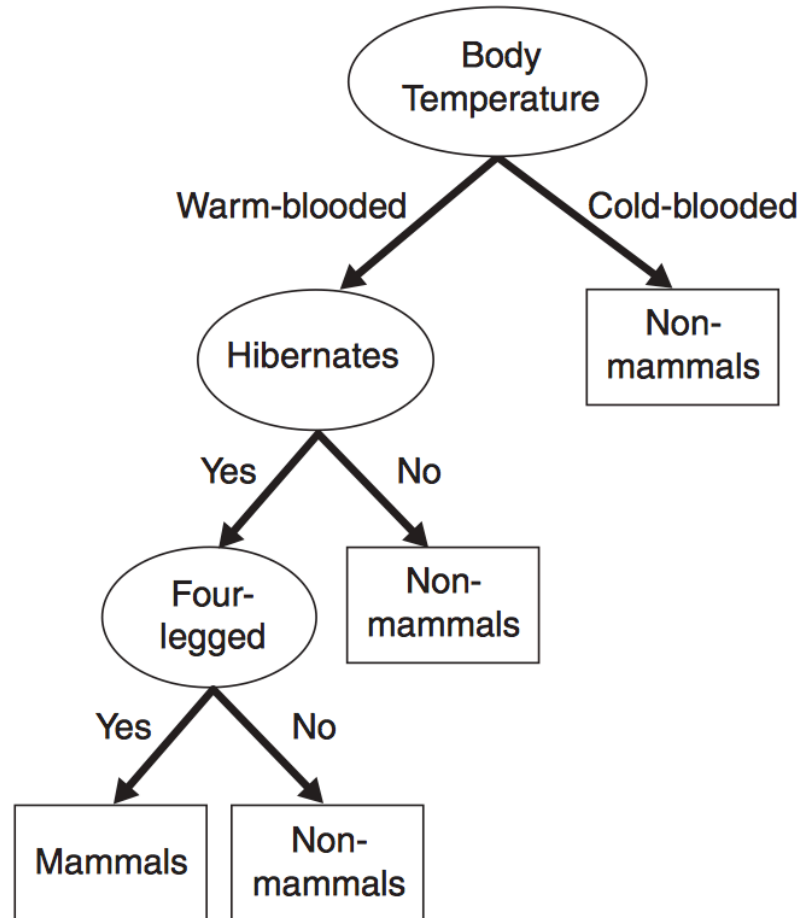
# Overfitting due to Insufficient Examples

**Table 4.5.** An example training set for classifying mammals.

| Name | Body Temperature | Gives Birth | Four-legged | Hibernates | Class Label |
|------|------------------|-------------|-------------|------------|-------------|
| salamander | cold-blooded | no | yes | yes | no |
| guppy | cold-blooded | yes | no | no | no |
| eagle | warm-blooded | no | no | no | no |
| poorwill | warm-blooded | no | no | yes | no |
| platypus | warm-blooded | no | yes | yes | yes |

Models that make their classification decisions based on a small number of training records are also susceptible to overfitting.

# Overfitting due to Insufficient Examples



**Figure 4.26.** Decision tree induced from the data set shown in Table 4.5.

# Overfitting due to Insufficient Examples

- All of these training records are labeled correctly and its training error is zero, its error rate on the test set is 30%.

- Humans, elephants, and dolphins are misclassified because the decision tree classifies all warm-blooded vertebrates that do not hibernate as non-mammals.

- The tree arrives at this classification decision because there is only one training record, which is an eagle, with such characteristics.

- This example clearly demonstrates the danger of making wrong predictions when there are not enough representative examples at the leaf nodes of a decision tree.

# Notes on Overfitting

- Overfitting results in decision trees that are more complex than necessary

- Training error no longer provides a good estimate of how well the tree will perform on previously unseen records

- Need new ways for estimating errors

# Estimating Generalization Errors

- Re-substitution errors: error on training ($\Sigma$ e(t) )
- Generalization errors: error on testing ($\Sigma$ e'(t))

- Methods for estimating generalization errors:
  - Optimistic approach:  e'(t) = e(t)
  - Pessimistic approach:
    - For each leaf node: e'(t) = (e(t)+0.5)
    - Total errors: e'(T) = e(T) + N $\times$ 0.5 (N: number of leaf nodes)
    - For a tree with 30 leaf nodes and 10 errors on training (out of 1000 instances):
      Training error = 10/1000 = 1%

      Generalization error = (10 + 30$\times$0.5)/1000 = 2.5%
  - Reduced error pruning (REP):
    - uses validation data set to estimate generalization error

# Occam's Razor

- Given two models of similar generalization errors, one should prefer the simpler model over the more complex model

- For complex models, there is a greater chance that it was fitted accidentally by errors in data

- Therefore, one should include model complexity when evaluating a model

# Handling Missing Attribute Values

- Missing values affect decision tree construction in three different ways:

  - Affects how impurity measures are computed

  - Affects how to distribute instance with missing value to child nodes

  - Affects how a test instance with missing value is classified

# Metrics for Performance Evaluation

- Focus on the predictive capability of a model
  - Rather than how fast it takes to classify or build models, scalability, etc.

- Confusion Matrix:

| | PREDICTED CLASS | | |
|---|---|---|---|
| ACTUAL CLASS | | Class=Yes | Class=No |
| | Class=Yes | a | b |
| | Class=No | c | d |

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

# Metrics for Performance Evaluation…

| | PREDICTED CLASS | | |
|---|---|---|---|
| **ACTUAL CLASS** | | Class=Yes | Class=No |
| | Class=Yes | a (TP) | b (FN) |
| | Class=No | c (FP) | d (TN) |

● Most widely-used metric:

$$\text{Accuracy} = \frac{a+d}{a+b+c+d} = \frac{TP+TN}{TP+TN+FP+FN}$$

# Limitation of Accuracy

- Consider a 2-class problem
    - Number of Class 0 examples = 9990
    - Number of Class 1 examples = 10

- If model predicts everything to be class 0, accuracy is 9990/10000 = 99.9 %
    - Accuracy is misleading because model does not detect any class 1 example

# Cost Matrix

| | PREDICTED CLASS | | |
|---|---|---|---|
| ACTUAL CLASS | C(i\|j) | **Class=Yes** | **Class=No** |
| | **Class=Yes** | C(Yes\|Yes) | C(No\|Yes) |
| | **Class=No** | C(Yes\|No) | C(No\|No) |

C(i|j): Cost of misclassifying class j example as class i

# Computing Cost of Classification

| Cost Matrix | PREDICTED CLASS | | |
|---|---|---|---|
| | C(i\|j) | + | - |
| ACTUAL CLASS | + | -1 | 100 |
| | - | 1 | 0 |

| Model M$_1$ | PREDICTED CLASS | | |
|---|---|---|---|
| | | + | - |
| ACTUAL CLASS | + | 150 | 40 |
| | - | 60 | 250 |

| Model M$_2$ | PREDICTED CLASS | | |
|---|---|---|---|
| | | + | - |
| ACTUAL CLASS | + | 250 | 45 |
| | - | 5 | 200 |

Accuracy = 80%

Cost = 3910

Accuracy = 90%

Cost = 4255

# Cost vs Accuracy

| Count | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| **ACTUAL CLASS** Class=Yes | a | b |
| Class=No | c | d |

Accuracy is proportional to cost if
1. $C(Yes|No)=C(No|Yes) = q$
2. $C(Yes|Yes)=C(No|No) = p$

$$N = a + b + c + d$$

$$Accuracy = (a + d)/N$$

| Cost | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| **ACTUAL CLASS** Class=Yes | p | q |
| Class=No | q | p |

$$Cost = p (a + d) + q (b + c)$$
$$= p (a + d) + q (N - a - d)$$
$$= q N - (q - p)(a + d)$$
$$= N [q - (q-p) \times Accuracy]$$

# Model Evaluation

- **Metrics for Performance Evaluation**
  - How to evaluate the performance of a model?

- <span style="color:red">**Methods for Performance Evaluation**</span>
  - How to obtain reliable estimates?

- **Methods for Model Comparison**
  - How to compare the relative performance among competing models?

# Methods for Performance Evaluation

- How to obtain a reliable estimate of performance?

- Performance of a model may depend on other factors besides the learning algorithm:
    - Class distribution
    - Cost of misclassification
    - Size of training and test sets

# Methods of Estimation

- Holdout
  - Reserve 2/3 for training and 1/3 for testing
- Random subsampling
  - Repeated holdout
- Cross validation
  - Partition data into k disjoint subsets
  - k-fold: train on k-1 partitions, test on the remaining one
  - Leave-one-out: k=n
- Stratified sampling
  - oversampling vs undersampling
- Bootstrap
  - Sampling with replacement

# Model Evaluation

- Metrics for Performance Evaluation
  - How to evaluate the performance of a model?

- Methods for Performance Evaluation
  - How to obtain reliable estimates?

- Methods for Model Comparison
  - How to compare the relative performance among competing models?

# Test of Significance

- Given two models:
  - Model M1: accuracy = 85%, tested on 30 instances
  - Model M2: accuracy = 75%, tested on 5000 instances

- Can we say M1 is better than M2?
  - How much confidence can we place on accuracy of M1 and M2?
  - Can the difference in performance measure be explained as a result of random fluctuations in the test set?

# Confidence Interval for Accuracy

- Prediction can be regarded as a Bernoulli trial
  - A Bernoulli trial has 2 possible outcomes
  - Possible outcomes for prediction: correct or wrong
  - Collection of Bernoulli trials has a Binomial distribution:
    - $x \sim Bin(N, p)$     x: number of correct predictions
    - e.g: Toss a fair coin 50 times, how many heads would turn up? Expected number of heads = $N \times p = 50 \times 0.5 = 25$

- Given x (# of correct predictions) or equivalently, acc=x/N, and N (# of test instances),

  Can we predict p (true accuracy of model)?

# Confidence Interval for Accuracy

- For large test sets (N > 30),
  - acc has a normal distribution with mean p and variance p(1-p)/N

$$P(Z_{\alpha/2} < \frac{acc - p}{\sqrt{p(1-p)/N}} < Z_{1-\alpha/2})$$

$$= 1 - \alpha$$

**Area = 1 - $\alpha$**

$Z_{\alpha/2}$  $Z_{1-\alpha/2}$

- Confidence Interval for p:

$$p = \frac{2 \times N \times acc + Z_{\alpha/2}^2 \pm \sqrt{Z_{\alpha/2}^2 + 4 \times N \times acc - 4 \times N \times acc^2}}{2(N + Z_{\alpha/2}^2)}$$

# Data Mining
# Cluster Analysis: Basic Concepts and Algorithms

Lecture Notes for Chapter 8

Introduction to Data Mining

by

Tan, Steinbach, Kumar

# Notion of a Cluster can be Ambiguous



How many clusters?

Six Clusters

Two Clusters

Four Clusters

# Types of Clusterings

- A clustering is a set of clusters

- Important distinction between hierarchical and partitional sets of clusters

- Partitional Clustering
  - A division data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset

- Hierarchical clustering
  - A set of nested clusters organized as a hierarchical tree

# Hierarchical Clustering



**Traditional Hierarchical Clustering**
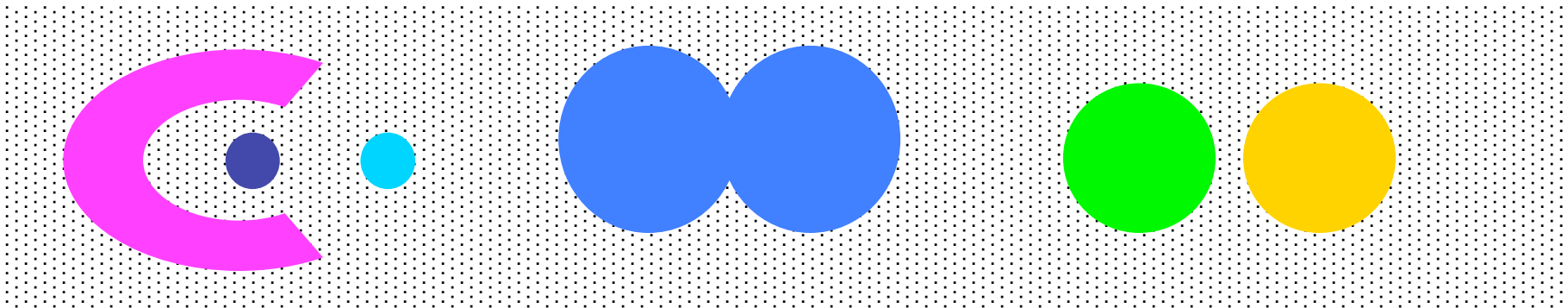
**Traditional Dendrogram**

**Non-traditional Hierarchical Clustering**

**Non-traditional Dendrogram**

# Types of Clusters: Density-Based

- ## Density-based

  - A cluster is a dense region of points, which is separated by low-density regions, from other regions of high density.

  - Used when the clusters are irregular or intertwined, and when noise and outliers are present.

**6 density-based clusters**

# Types of Clusters: Objective Function

- **Clusters Defined by an Objective Function**
    - Finds clusters that minimize or maximize an objective function.
    - Enumerate all possible ways of dividing the points into clusters and evaluate the `goodness' of each potential set of clusters by using the given objective function.  (NP Hard)
    - Can have global or local objectives.
        - Hierarchical clustering algorithms typically have local objectives
        - Partitional algorithms typically have global objectives
    - A variation of the global objective function approach is to fit the data to a parameterized model.
        - Parameters for the model are determined from the data.
        - Mixture models assume that the data is a 'mixture' of a number of statistical distributions.

# Types of Clusters: Objective Function …

- Map the clustering problem to a different domain and solve a related problem in that domain

    – Proximity matrix defines a weighted graph, where the nodes are the points being clustered, and the weighted edges represent the proximities between points

    – Clustering is equivalent to breaking the graph into connected components, one for each cluster.
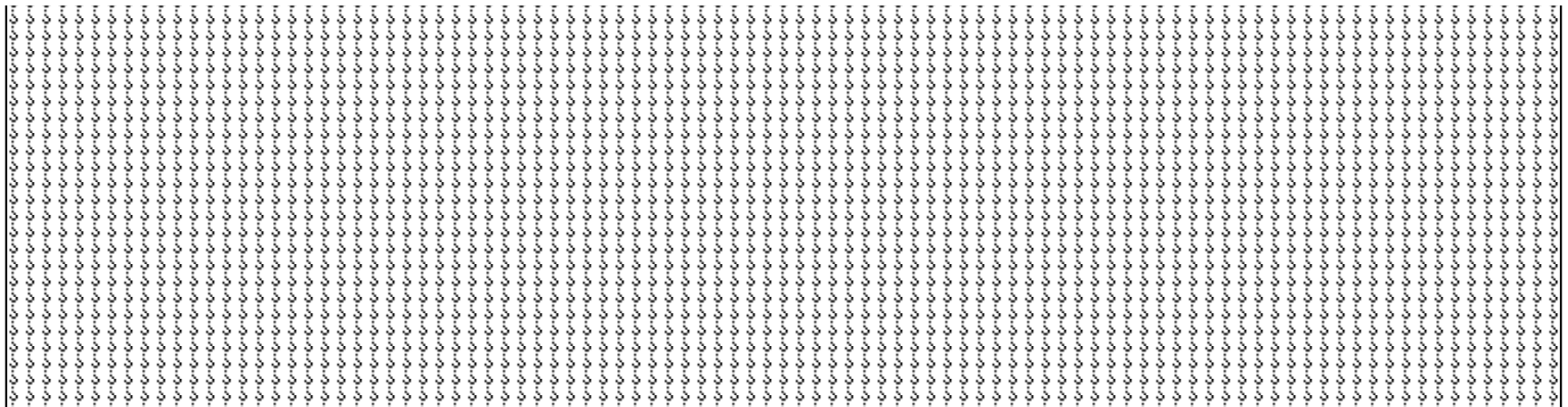
    – Want to minimize the edge weight between clusters and maximize the edge weight within

# Clustering Algorithms

- K-means and its variants

- Hierarchical clustering
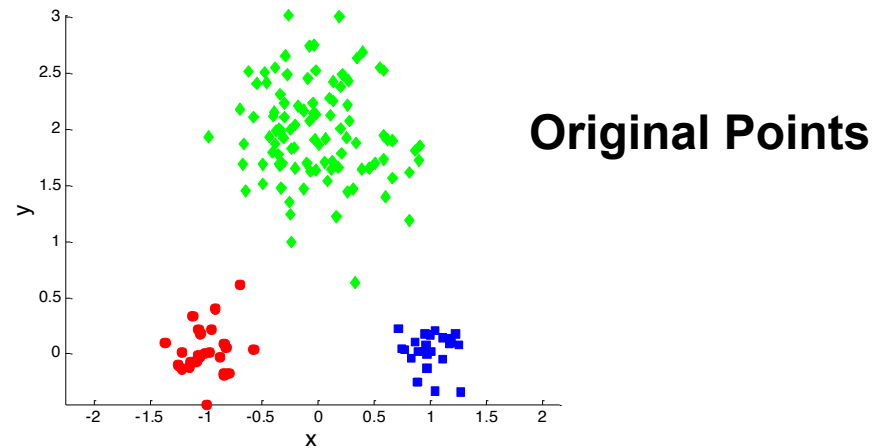
- Density-based clustering

# K-means Clustering

- Partitional clustering approach

- Each cluster is associated with a centroid (center point)

- Each point is assigned to the cluster with the closest centroid

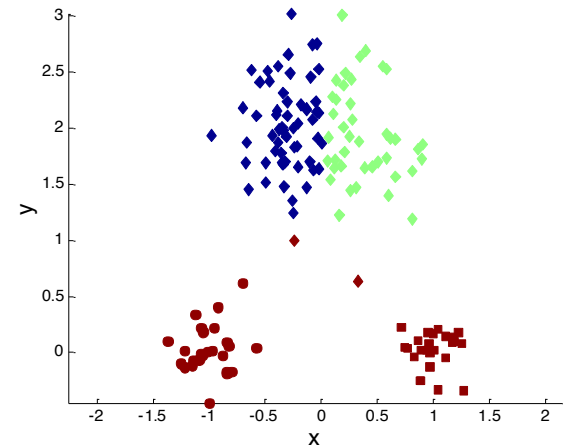- Number of clusters, K, must be specified

- The basic algorithm is very simple

# K-means Clustering – Details

- Initial centroids are often chosen randomly.
  - Clusters produced vary from one run to another.
- The centroid is (typically) the mean of the points in the cluster.
- 'Closeness' is measured by Euclidean distance, cosine similarity, correlation, etc.
- K-means will converge for common similarity measures mentioned above.
- Most of the convergence happens in the first few iterations.
  - Often the stopping condition is changed to 'Until relatively few points change clusters'
- Complexity is O( n * K * I * d )
  - n = number of points, K = number of clusters,
    I = number of iterations, d = number of attributes
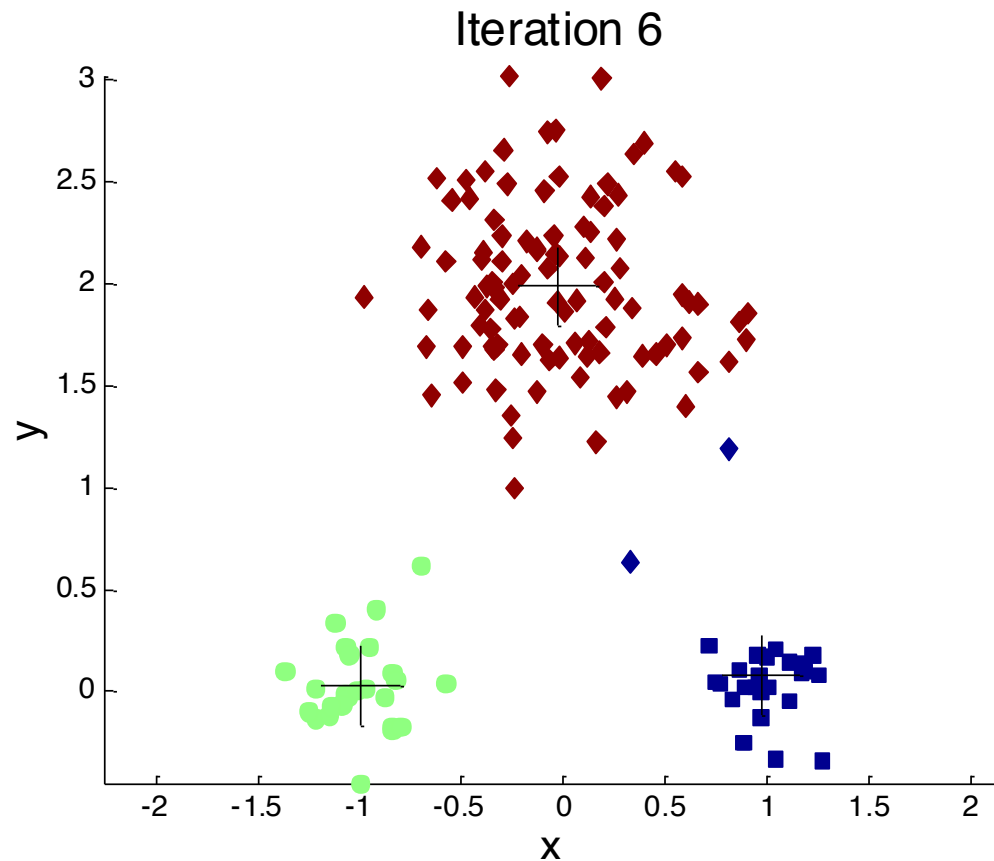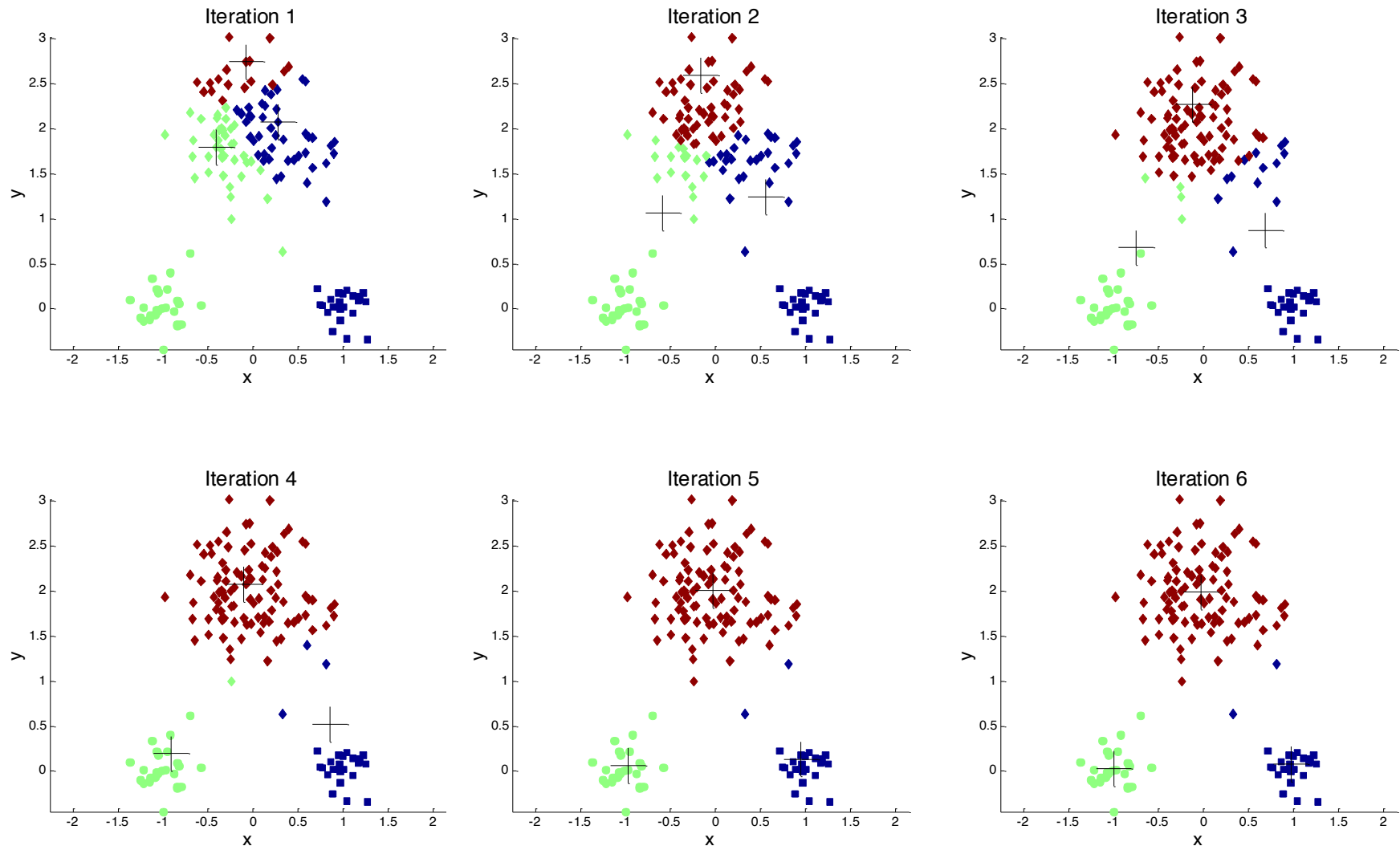
# Two different K-means Clusterings



**Original Points**

**Optimal Clustering**

**Sub-optimal Clustering**

# Importance of Choosing Initial Centroids



Iteration 6

# Importance of Choosing Initial Centroids
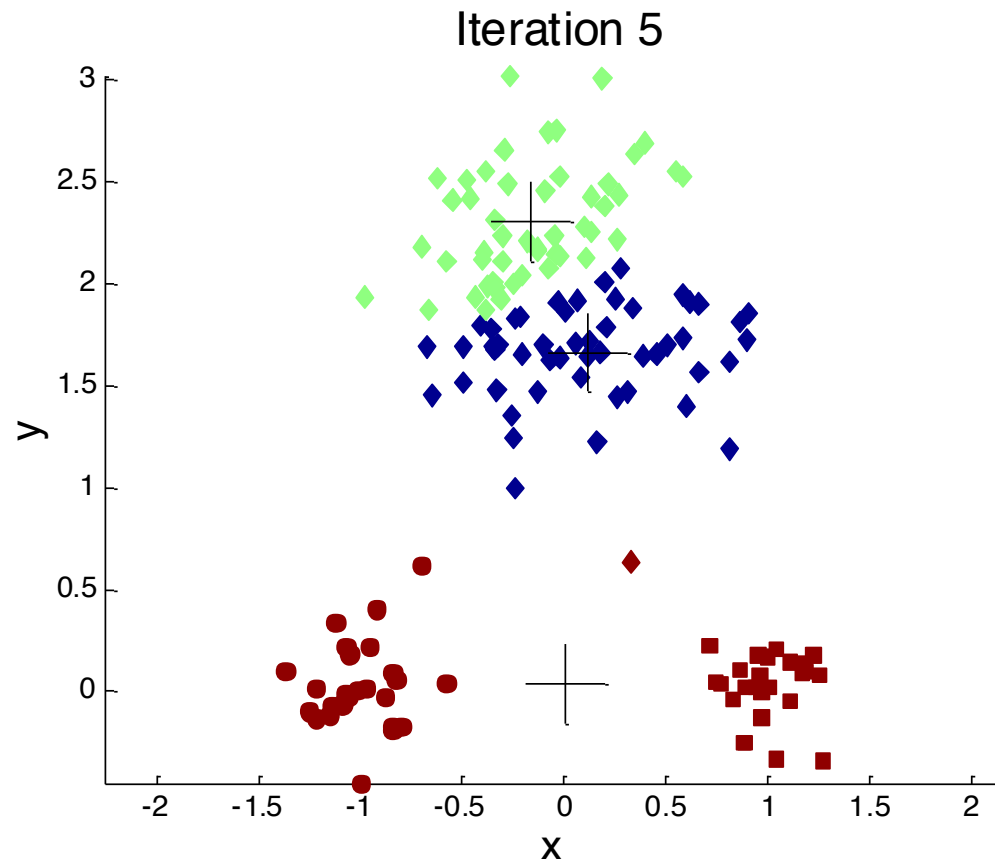
# Evaluating K-means Clusters

- Most common measure is Sum of Squared Error (SSE)
  - For each point, the error is the distance to the nearest cluster
  - To get SSE, we square these errors and sum them.
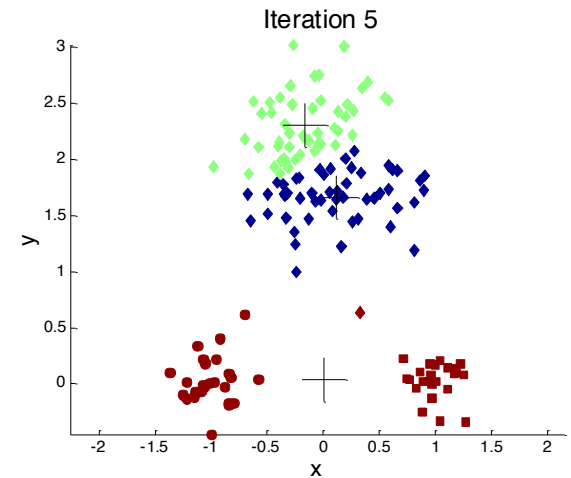
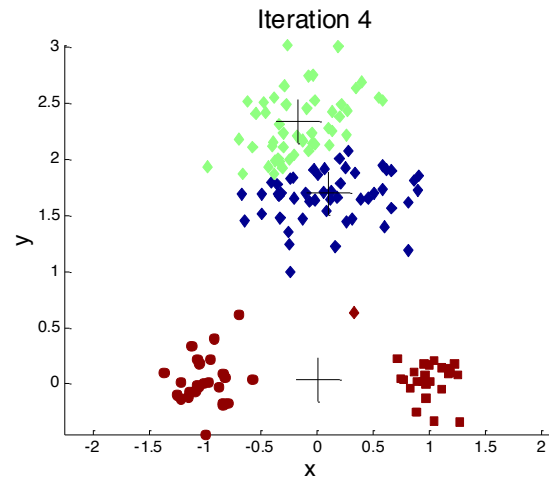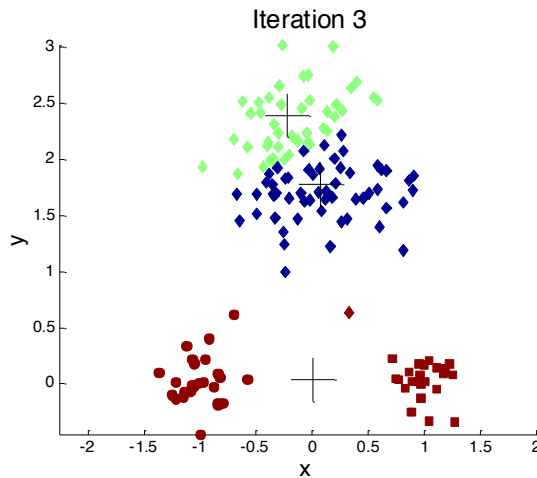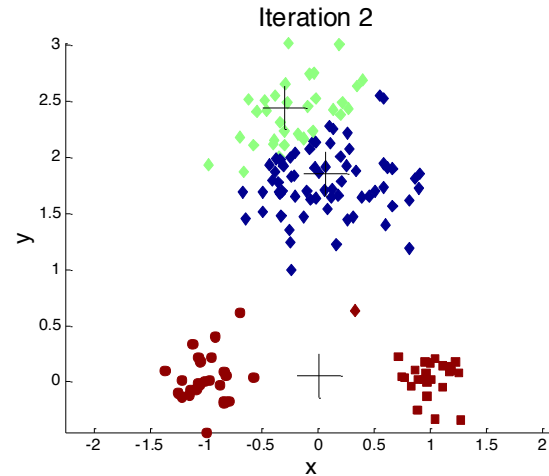$$SSE = \sum_{i=1}^{K} \sum_{x \in C_i} dist^2(m_i, x)$$

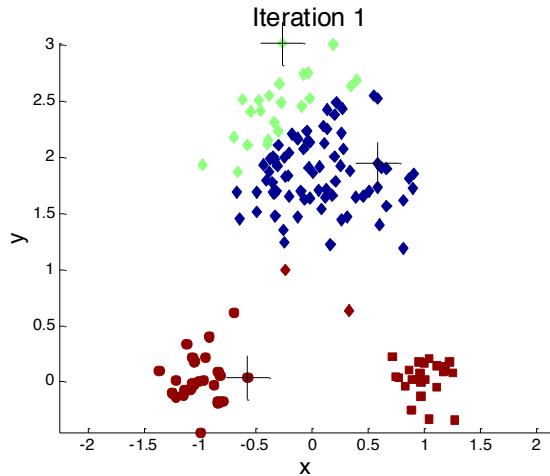  - $x$ is a data point in cluster $C_i$ and $m_i$ is the representative point for cluster $C_i$
    - can show that $m_i$ corresponds to the center (mean) of the cluster
  - Given two clusters, we can choose the one with the smallest error
  - One easy way to reduce SSE is to increase K, the number of clusters
    - A good clustering with smaller K can have a lower SSE than a poor clustering with higher K

# Importance of Choosing Initial Centroids ...



Iteration 5

# Importance of Choosing Initial Centroids …

# Solutions to Initial Centroids Problem

- Multiple runs
  - Helps, but probability is not on your side
- Sample and use hierarchical clustering to determine initial centroids
- Select more than k initial centroids and then select among these initial centroids
  - Select most widely separated
- Postprocessing
- Bisecting K-means
  - Not as susceptible to initialization issues
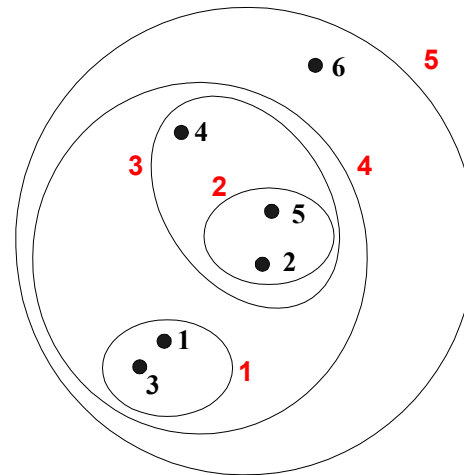
# Pre-processing and Post-processing

- Pre-processing
  - Normalize the data
  - Eliminate outliers

- Post-processing
  - Eliminate small clusters that may represent outliers
  - Split 'loose' clusters, i.e., clusters with relatively high SSE
  - Merge clusters that are 'close' and that have relatively low SSE
  - Can use these steps during the clustering process

# Hierarchical Clustering

- Produces a set of nested clusters organized as a hierarchical tree

- Can be visualized as a dendrogram
  - A tree like diagram that records the sequences of merges or splits

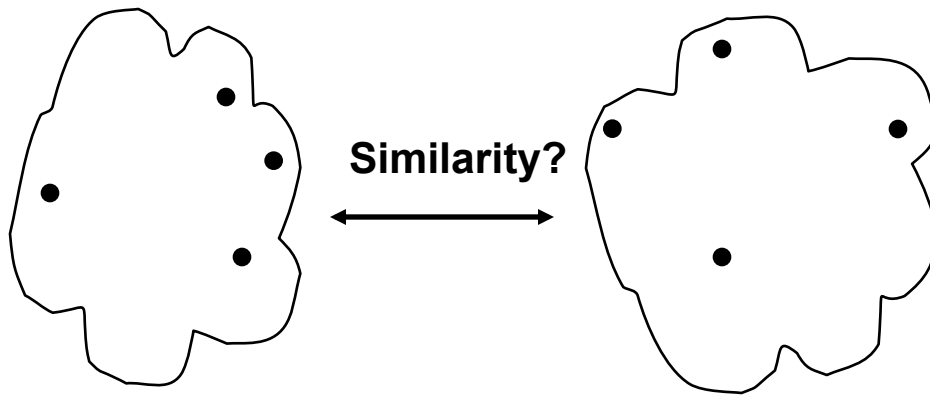# Strengths of Hierarchical Clustering

- Do not have to assume any particular number of clusters
  - Any desired number of clusters can be obtained by 'cutting' the dendogram at the proper level

- They may correspond to meaningful taxonomies
  - Example in biological sciences (e.g., animal kingdom, phylogeny reconstruction, …)

# Hierarchical Clustering

- Two main types of hierarchical clustering
    - Agglomerative:
        - Start with the points as individual clusters
        - At each step, merge the closest pair of clusters until only one cluster (or k clusters) left

    - Divisive:
        - Start with one, all-inclusive cluster
        - At each step, split a cluster until each cluster contains a point (or there are k clusters)

- Traditional hierarchical algorithms use a similarity or distance matrix
    - Merge or split one cluster at a time

# How to Define Inter-Cluster Similarity



Similarity?

| | p1 | p2 | p3 | p4 | p5 | . . . |
|---|---|---|---|---|---|---|
| **p1** | | | | | | |
| **p2** | | | | | | |
| **p3** | | | | | | |
| **p4** | | | | | | |
| **p5** | | | | | | |
| **.** | | | | | | |
| **.** | | | | | | |
| **.** | | | | | | |

**Proximity Matrix**

- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



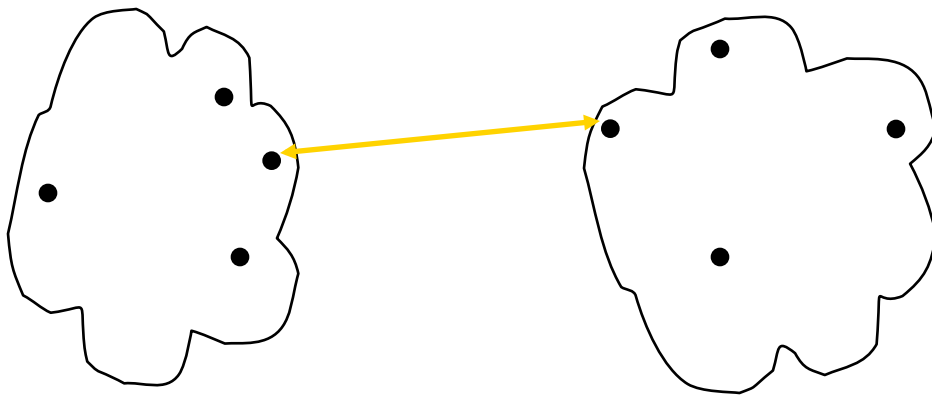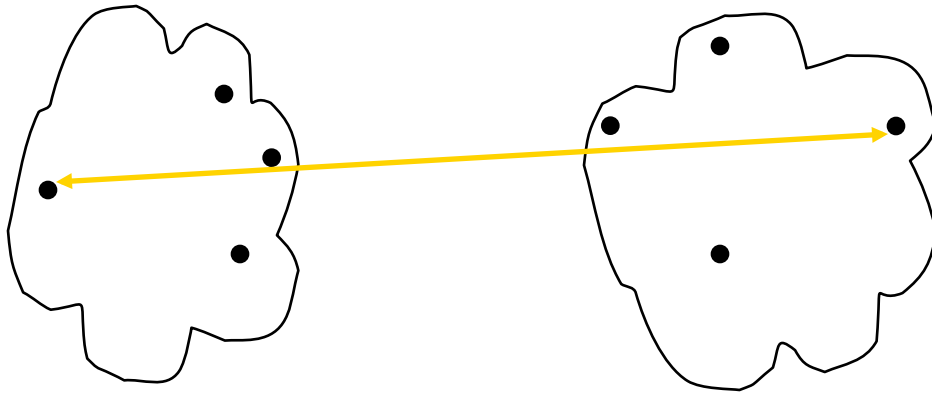|    | p1 | p2 | p3 | p4 | p5 | . . . |
|----|----|----|----|----|----|-------|
| p1 |    |    |    |    |    |       |
| p2 |    |    |    |    |    |       |
| p3 |    |    |    |    |    |       |
| p4 |    |    |    |    |    |       |
| p5 |    |    |    |    |    |       |
| .  |    |    |    |    |    |       |

**Proximity Matrix**

- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity

| | p1 | p2 | p3 | p4 | p5 | . . . |
|-----|-----|-----|-----|-----|-----|-----|
| p1 | | | | | | |
| p2 | | | | | | |
| p3 | | | | | | |
| p4 | | | | | | |
| p5 | | | | | | |
| . | | | | | | |

**Proximity Matrix**
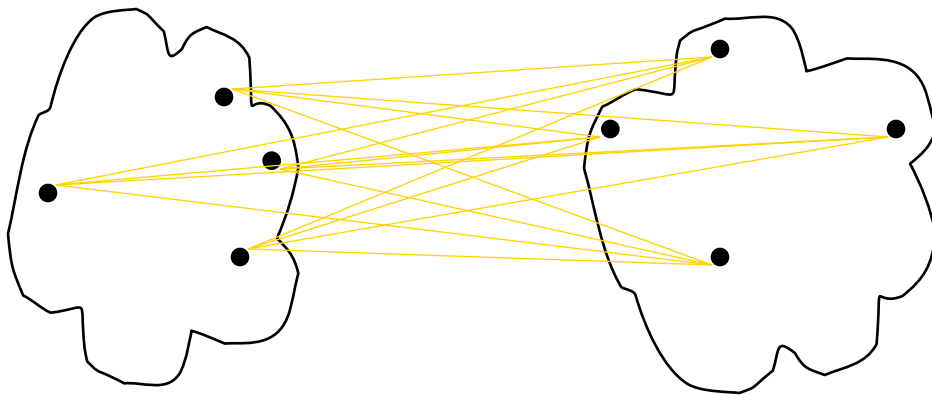
- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



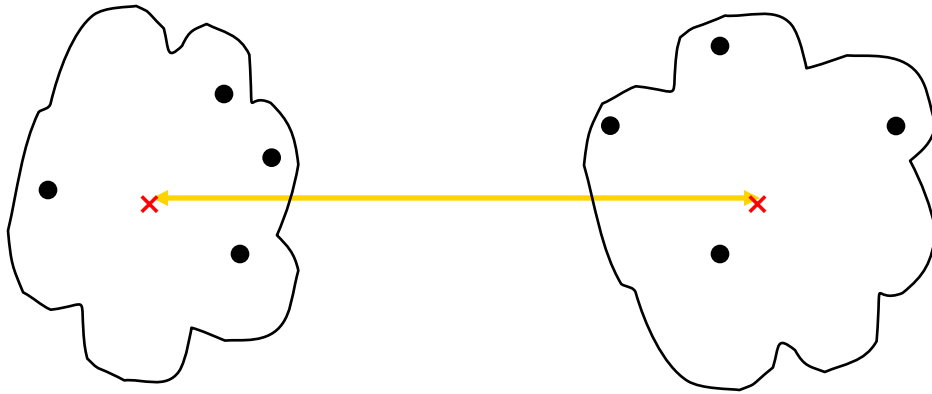|    | p1 | p2 | p3 | p4 | p5 | . . . |
|----|----|----|----|----|----|-------|
| p1 |    |    |    |    |    |       |
| p2 |    |    |    |    |    |       |
| p3 |    |    |    |    |    |       |
| p4 |    |    |    |    |    |       |
| p5 |    |    |    |    |    |       |
| .  |    |    |    |    |    |       |

**Proximity Matrix**

- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
  - Ward's Method uses squared error

# How to Define Inter-Cluster Similarity



|    | p1 | p2 | p3 | p4 | p5 | . . . |
|----|----|----|----|----|----|-------|
| p1 |    |    |    |    |    |       |
| p2 |    |    |    |    |    |       |
| p3 |    |    |    |    |    |       |
| p4 |    |    |    |    |    |       |
| p5 |    |    |    |    |    |       |
| .  |    |    |    |    |    |       |
| .  |    |    |    |    |    |       |
| .  |    |    |    |    |    |       |

**Proximity Matrix**
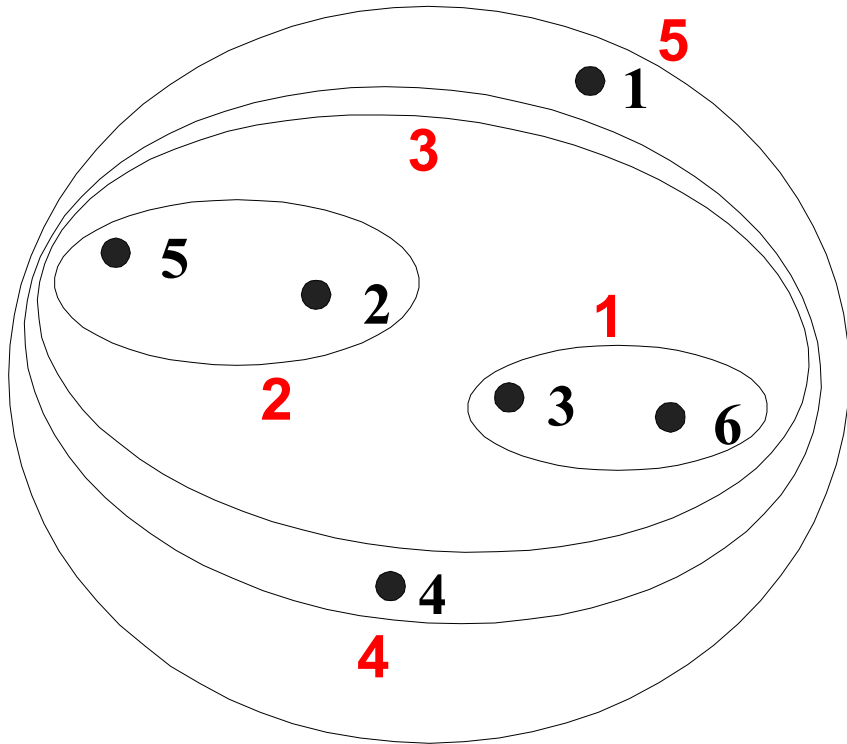
- MIN
- MAX
- Group Average
- Distance Between Centroids
- Other methods driven by an objective function
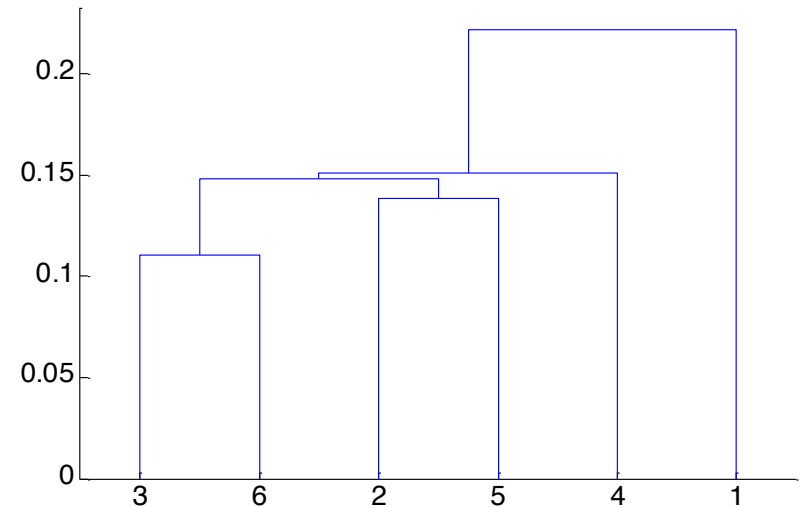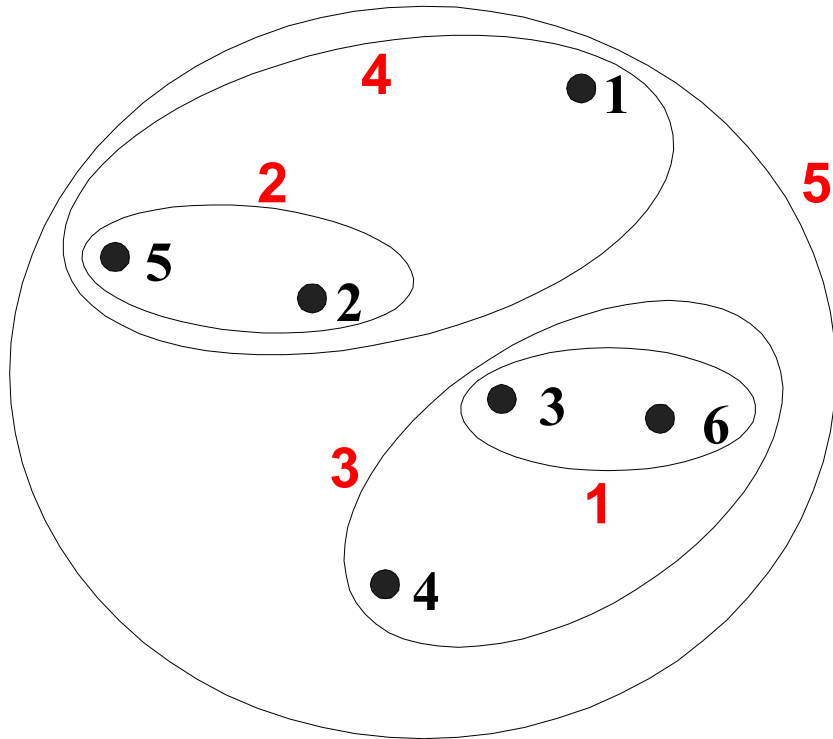    - Ward's Method uses squared error
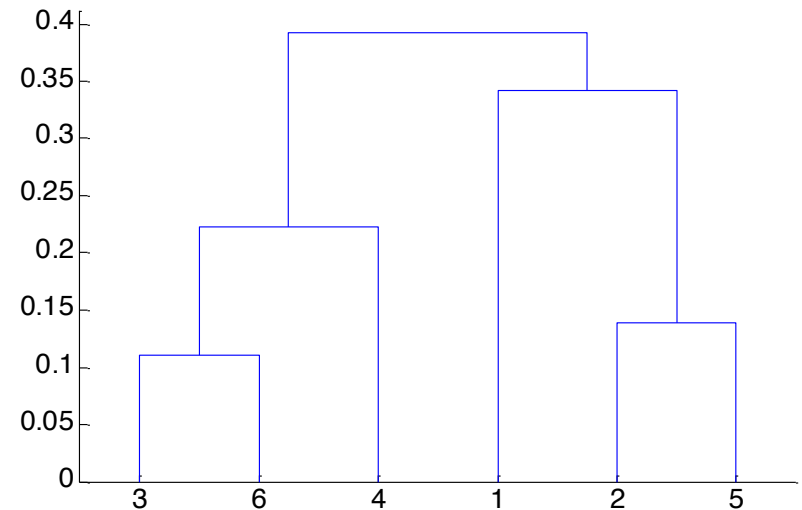
# Hierarchical Clustering: MIN



**Nested Clusters**

**Dendrogram**

# Hierarchical Clustering: MAX



**Nested Clusters**

**Dendrogram**

# Cluster Validity

- For supervised classification we have a variety of measures to evaluate how good our model is
    – Accuracy, precision, recall

- For cluster analysis, the analogous question is how to evaluate the "goodness" of the resulting clusters?

- But "clusters are in the eye of the beholder"!

- Then why do we want to evaluate them?
    – To avoid finding patterns in noise
    – To compare clustering algorithms
    – To compare two sets of clusters
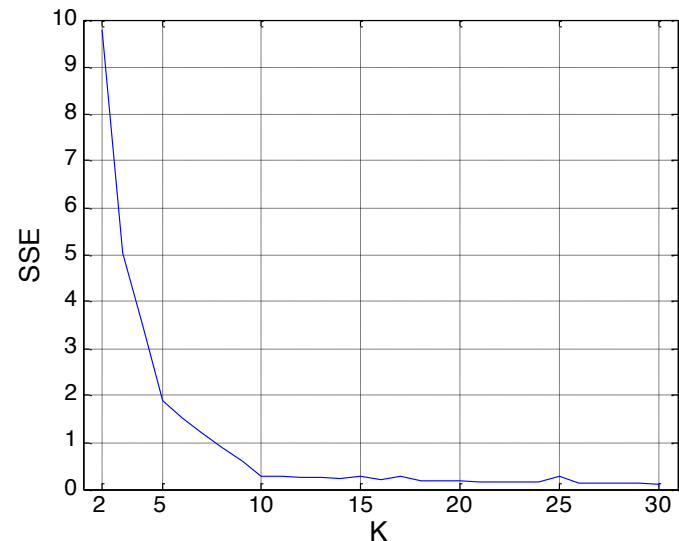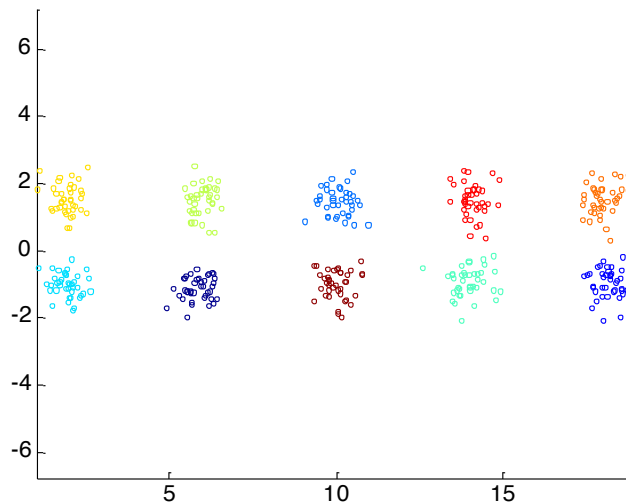    – To compare two clusters

# Measures of Cluster Validity

● Numerical measures that are applied to judge various aspects of cluster validity, are classified into the following three types.

– **External Index:** Used to measure the extent to which cluster labels match externally supplied class labels.
  ◆ Entropy

– **Internal Index:** Used to measure the goodness of a clustering structure *without* respect to external information.
  ◆ Sum of Squared Error (SSE)

– **Relative Index:** Used to compare two different clusterings or clusters.
  ◆ Often an external or internal index is used for this function, e.g., SSE or entropy

● Sometimes these are referred to as criteria instead of indices

– However, sometimes criterion is the general strategy and index is the numerical measure that implements the criterion.
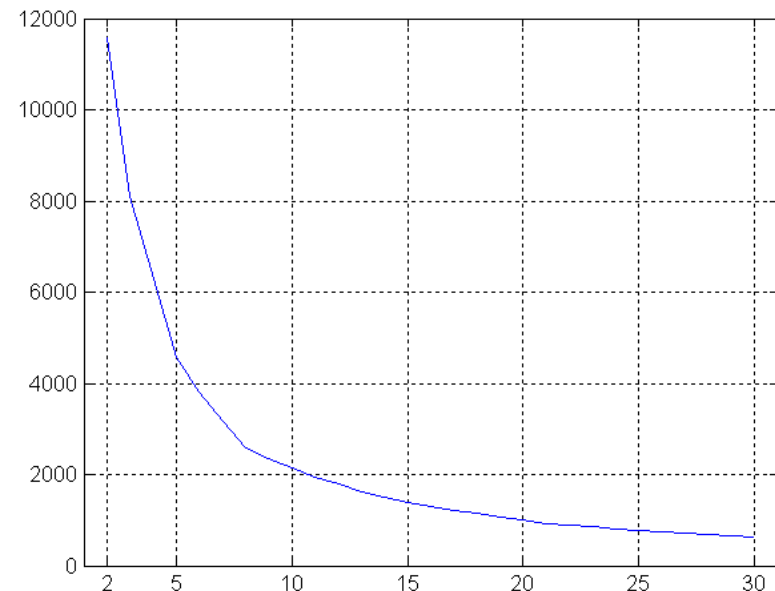
# Internal Measures: SSE

- Clusters in more complicated figures aren't well separated
- Internal Index:  Used to measure the goodness of a clustering structure without respect to external information
  - SSE
- SSE is good for comparing two clusterings or two clusters (average SSE).
- Can also be used to estimate the number of clusters

# Internal Measures: SSE

- SSE curve for a more complicated data set



**SSE of clusters found using K-means**

# Internal Measures: Cohesion and Separation

- **Cluster Cohesion**: Measures how closely related are objects in a cluster
  - Example: SSE
- **Cluster Separation**: Measure how distinct or well-separated a cluster is from other clusters
- Example: Squared Error
  - Cohesion is measured by the within cluster sum of squares (SSE)

$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$

  - Separation is measured by the between cluster sum of squares

$$BSS = \sum_i |C_i| (m - m_i)^2$$

   - Where $|C_i|$ is the size of cluster i

# Internal Measures: Cohesion and Separation

- Example: SSE
  - BSS + WSS = constant



**K=1 cluster:**

$$WSS = (1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2 = 10$$

$$BSS = 4 \times (3-3)^2 = 0$$

$$Total = 10 + 0 = 10$$

**K=2 clusters:**

$$WSS = (1-1.5)^2 + (2-1.5)^2 + (4-4.5)^2 + (5-4.5)^2 = 1$$

$$BSS = 2 \times (3-1.5)^2 + 2 \times (4.5-3)^2 = 9$$

$$Total = 1 + 9 = 10$$

# Internal Measures: Cohesion and Separation

- A proximity graph based approach can also be used for cohesion and separation.
    - Cluster cohesion is the sum of the weight of all links within a cluster.
    - Cluster separation is the sum of the weights between nodes in the cluster and nodes outside the cluster.

cohesion                                    separation

# Internal Measures: Silhouette Coefficient

- Silhouette Coefficient combine ideas of both cohesion and separation, but for individual points, as well as clusters and clusterings
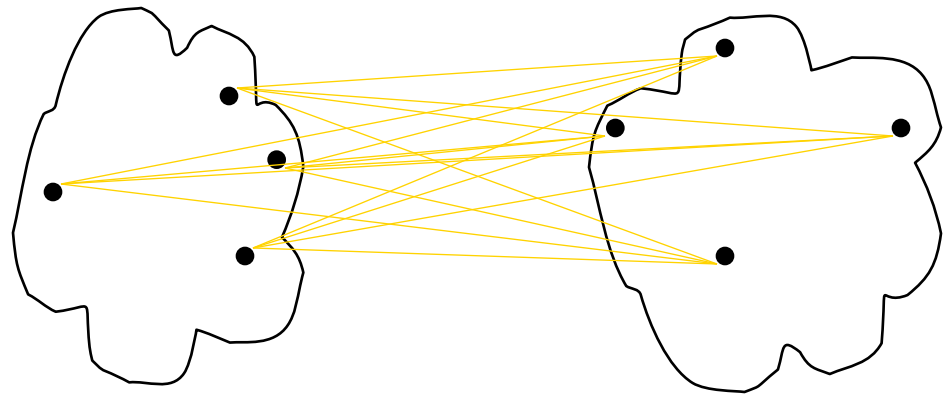- For an individual point, $i$
    - Calculate $a$ = average distance of $i$ to the points in its cluster
    - Calculate $b$ = min (average distance of $i$ to points in another cluster)
    - The silhouette coefficient for a point is then given by

        $s = 1 - a/b$   if $a < b$,   (or $s = b/a - 1$   if $a \geq b$, not the usual case)

    - Typically between 0 and 1.
    - The closer to 1 the better.

- Can calculate the Average Silhouette width for a cluster or a clustering

# Definition: Frequent Itemset

- **Itemset**
    - A collection of one or more items
        - Example: {Milk, Bread, Diaper}
    - k-itemset
        - An itemset that contains k items
- **Support count ($\sigma$)**
    - Frequency of occurrence of an itemset
    - E.g. $\sigma(\{Milk, Bread, Diaper\}) = 2$
- **Support**
    - Fraction of transactions that contain an itemset
    - E.g. s({Milk, Bread, Diaper}) = 2/5
- **Frequent Itemset**
    - An itemset whose support is greater than or equal to a *minsup* threshold

| TID | Items |
|-----|-------|
| 1 | **Bread, Milk** |
| 2 | **Bread, Diaper, Beer, Eggs** |
| 3 | **Milk, Diaper, Beer, Coke** |
| 4 | **Bread, Milk, Diaper, Beer** |
| 5 | **Bread, Milk, Diaper, Coke** |

# Definition: Association Rule

- **Association Rule**
  - An implication expression of the form $X \rightarrow Y$, where X and Y are itemsets
  - Example:
    {Milk, Diaper} $\rightarrow$ {Beer}

- **Rule Evaluation Metrics**
  - Support (s)
    - Fraction of transactions that contain both X and Y
  - Confidence (c)
    - Measures how often items in Y appear in transactions that contain X

| TID | Items |
|---|---|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

Example:

$$\{Milk, Diaper\} \Rightarrow Beer$$

$$s = \frac{\sigma(Milk, Diaper, Beer)}{|T|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(Milk, Diaper, Beer)}{\sigma(Milk, Diaper)} = \frac{2}{3} = 0.67$$

# Mining Association Rules

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

## Example of Rules:

{Milk,Diaper} $\rightarrow$ {Beer} (s=0.4, c=0.67)
{Milk,Beer} $\rightarrow$ {Diaper} (s=0.4, c=1.0)
{Diaper,Beer} $\rightarrow$ {Milk} (s=0.4, c=0.67)
{Beer} $\rightarrow$ {Milk,Diaper} (s=0.4, c=0.67)
{Diaper} $\rightarrow$ {Milk,Beer} (s=0.4, c=0.5)
{Milk} $\rightarrow$ {Diaper,Beer} (s=0.4, c=0.5)

## Observations:
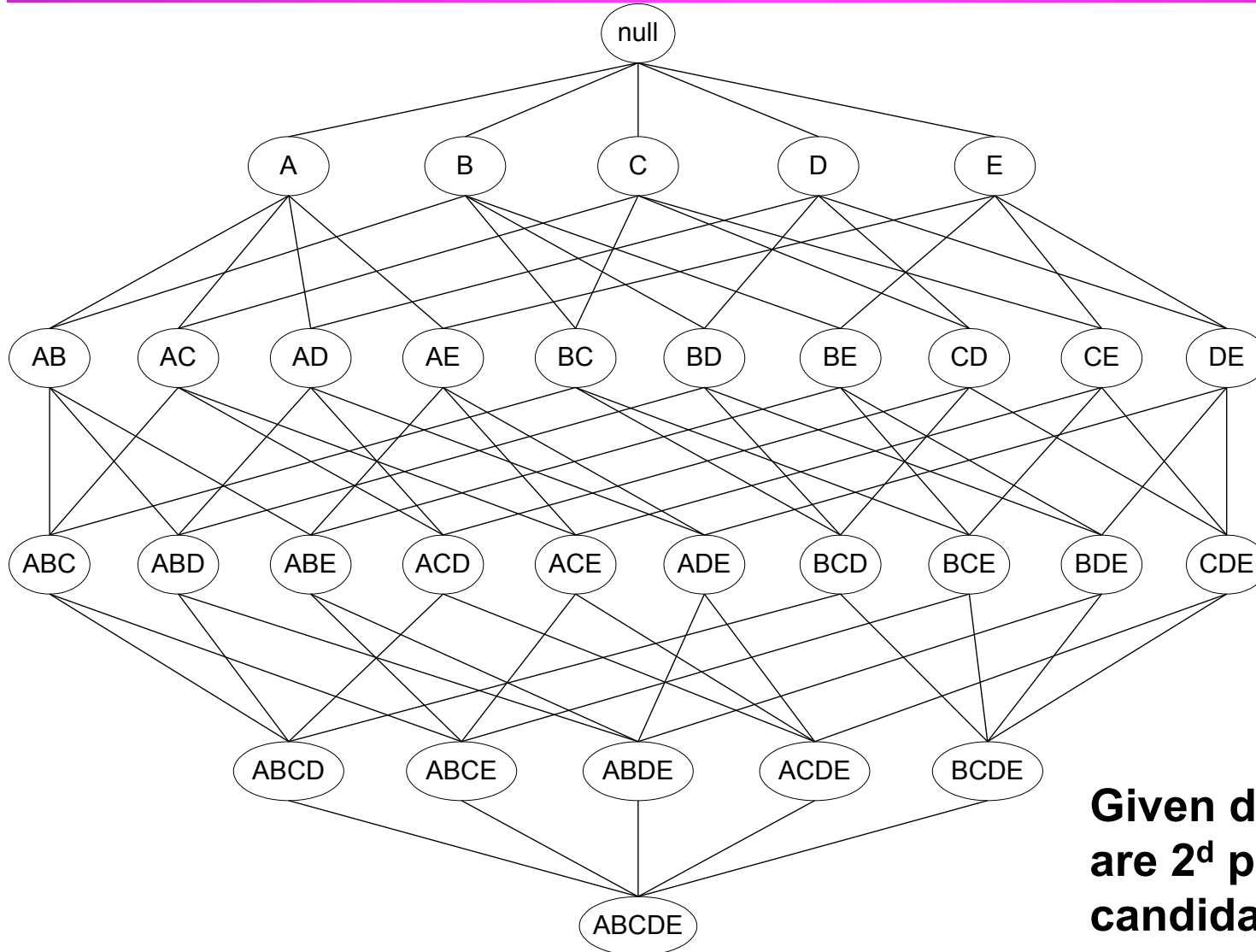
• All the above rules are binary partitions of the same itemset:
  {Milk, Diaper, Beer}

• Rules originating from the same itemset have identical support but can have different confidence

• Thus, we may decouple the support and confidence requirements

# Mining Association Rules

- Two-step approach:

  1. **Frequent Itemset Generation**
     - Generate all itemsets whose support $\geq$ minsup

  2. **Rule Generation**
     - Generate high confidence rules from each frequent itemset, where each rule is a binary partitioning of a frequent itemset

- Frequent itemset generation is still computationally expensive

# Frequent Itemset Generation



**Given d items, there are $2^d$ possible candidate itemsets**

# Frequent Itemset Generation

- Brute-force approach:
  - Each itemset in the lattice is a candidate frequent itemset
  - Count the support of each candidate by scanning the database

| TID | Items |
|-----|-------|
| 1 | Bread, Milk |
| 2 | Bread, Diaper, Beer, Eggs |
| 3 | Milk, Diaper, Beer, Coke |
| 4 | Bread, Milk, Diaper, Beer |
| 5 | Bread, Milk, Diaper, Coke |

N

w

M

- Match each transaction against every candidate

# Computational Complexity

● Given d unique items:
- Total number of itemsets = $2^d$
- Total number of possible association rules:



$$R = \sum_{k=1}^{d-1}\left[\binom{d}{k} \times \sum_{j=1}^{d-k}\binom{d-k}{j}\right]$$

$$= 3^d - 2^{d+1} + 1$$

**If d=6, R = 602 rules**

# Frequent Itemset Generation Strategies

- Reduce the number of candidates (M)
  - Complete search: $M = 2^d$
  - Use pruning techniques to reduce M

- Reduce the number of transactions (N)
  - Reduce size of N as the size of itemset increases
  - Used by DHP and vertical-based mining algorithms

- Reduce the number of comparisons (NM)
  - Use efficient data structures to store the candidates or transactions
  - No need to match every candidate against

# Reducing Number of Candidates

- **Apriori principle**:
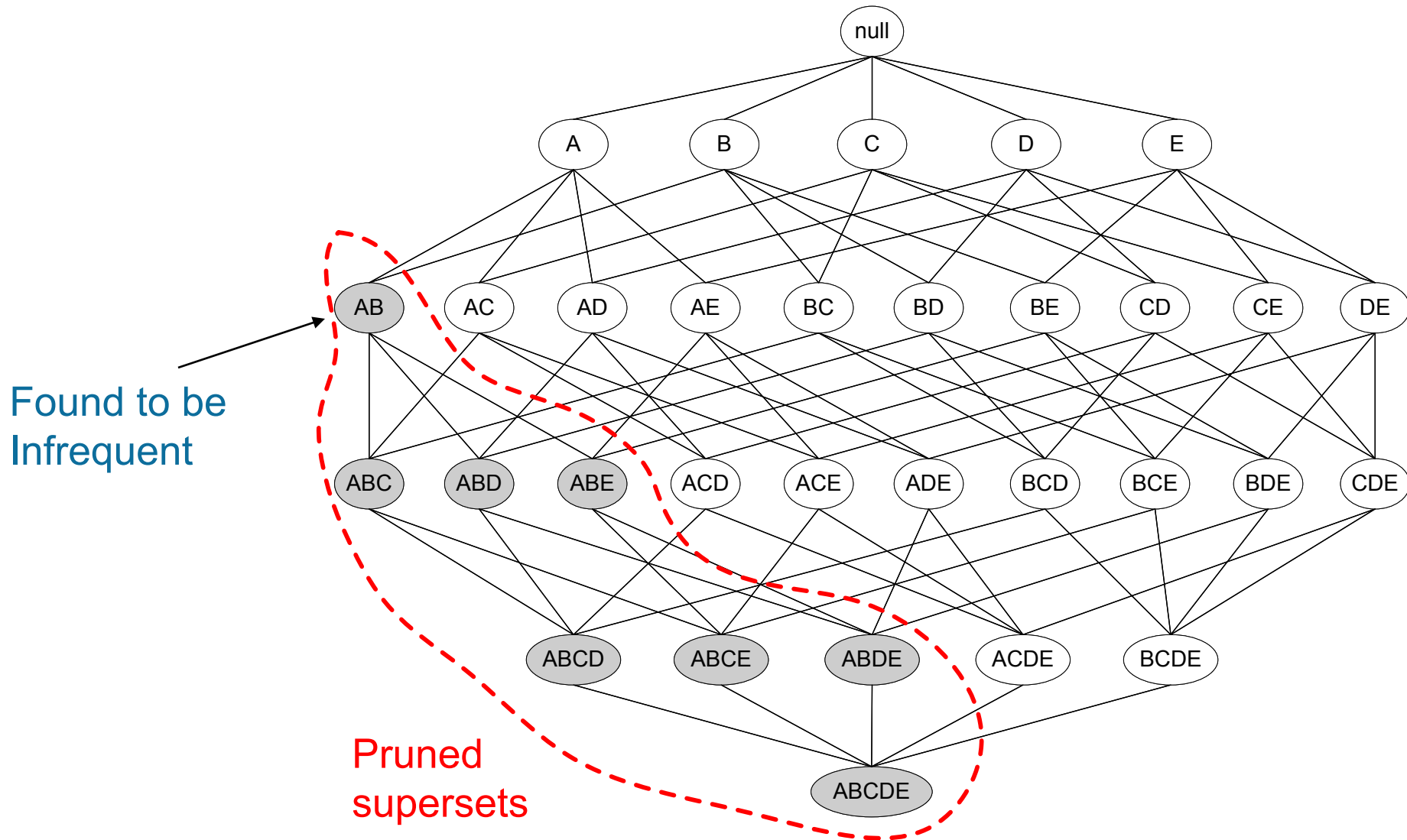  - If an itemset is frequent, then all of its subsets must also be frequent

- Apriori principle holds due to the following property of the support measure:

$$\forall X, Y : (X \subseteq Y) \Rightarrow s(X) \geq s(Y)$$

  - Support of an itemset never exceeds the support of its subsets
  - This is known as the anti-monotone property of support

# Illustrating Apriori Principle



Found to be Infrequent

Pruned supersets

# Illustrating Apriori Principle

| Item | Count |
|------|-------|
| **Bread** | 4 |
| Coke | 2 |
| **Milk** | 4 |
| **Beer** | 3 |
| **Diaper** | 4 |
| Eggs | 1 |

Items (1-itemsets)

| Itemset | Count |
|---------|-------|
| **{Bread,Milk}** | 3 |
| **{Bread,Beer}** | 2 |
| **{Bread,Diaper}** | 3 |
| **{Milk,Beer}** | 2 |
| **{Milk,Diaper}** | 3 |
| **{Beer,Diaper}** | 3 |

Pairs (2-itemsets)

(No need to generate candidates involving Coke or Eggs)

**Minimum Support = 3**

Triplets (3-itemsets)

| Itemset | Count |
|---------|-------|
| **{Bread,Milk,Diaper}** | 3 |

If every subset is considered,
$$^6C_1 + {}^6C_2 + {}^6C_3 = 41$$
With support-based pruning,
$$6 + 6 + 1 = 13$$

# Apriori Algorithm

● Method:

– Let k=1

– Generate frequent itemsets of length 1

– Repeat until no new frequent itemsets are identified

◆ Generate length (k+1) candidate itemsets from length k frequent itemsets

◆ Prune candidate itemsets containing subsets of length k that are infrequent

◆ Count the support of each candidate by scanning the DB

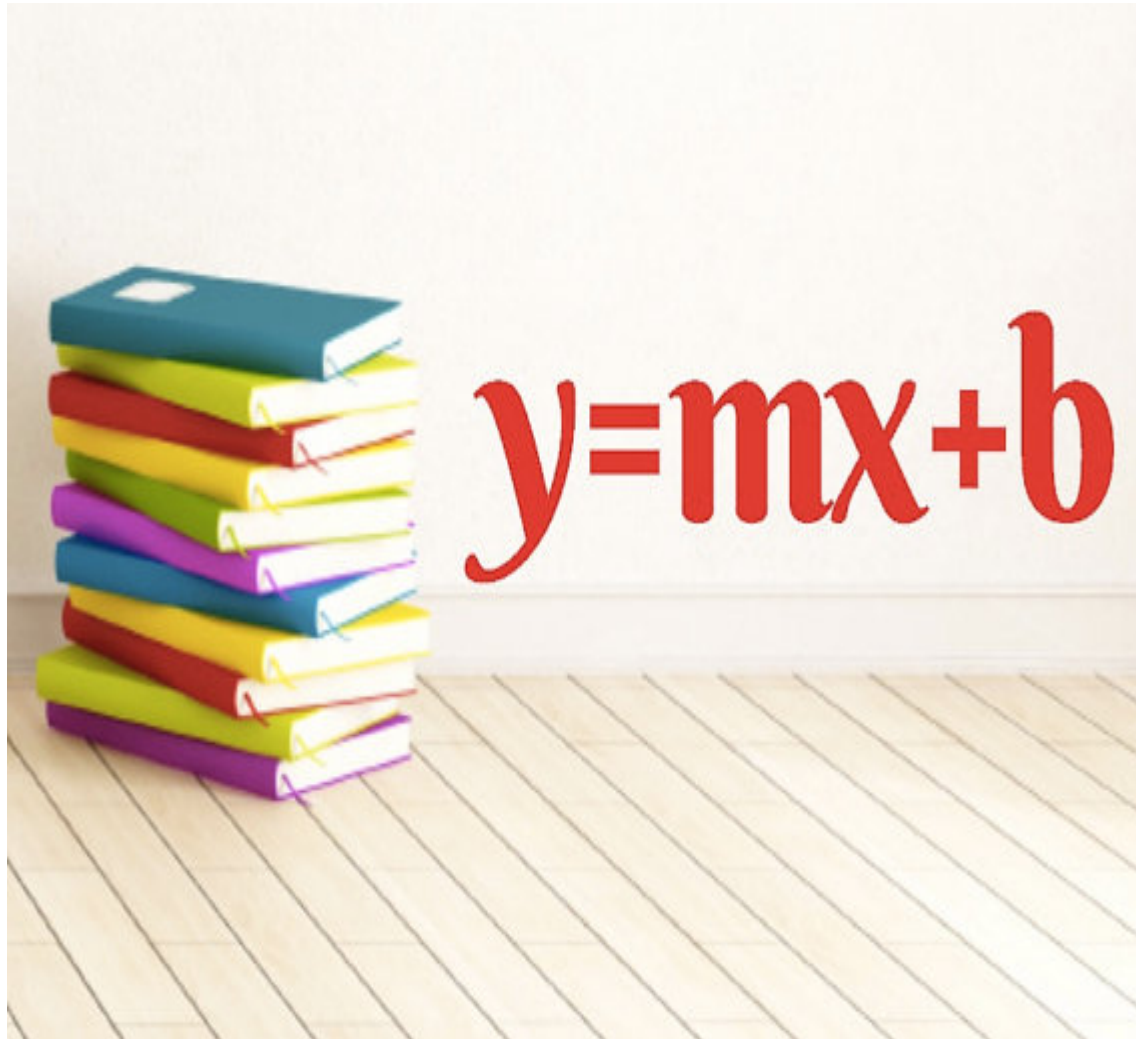◆ Eliminate candidates that are infrequent, leaving only those that are frequent

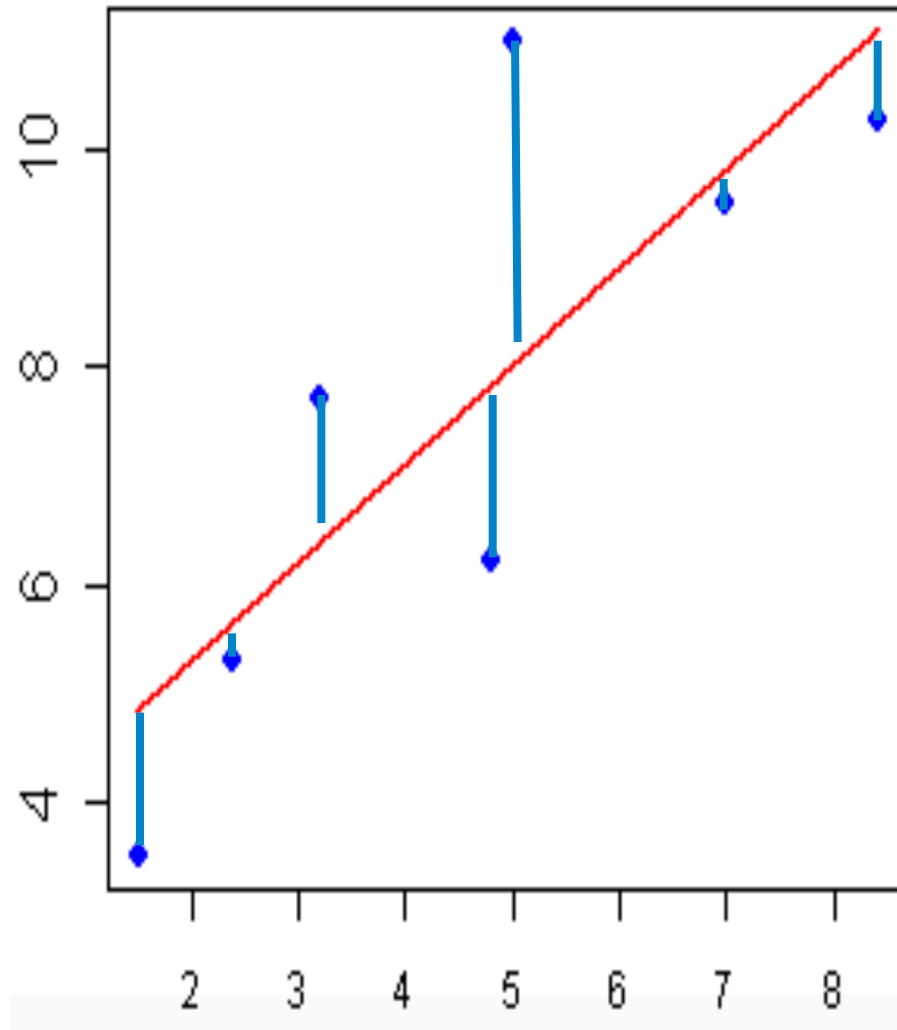# Challenges of Data Mining

- Scalability

- Dimensionality

- Complex and Heterogeneous Data

- Data Quality

- Data Ownership and Distribution

- Privacy Preservation

- Streaming Data

# Linear regression

# Residuals

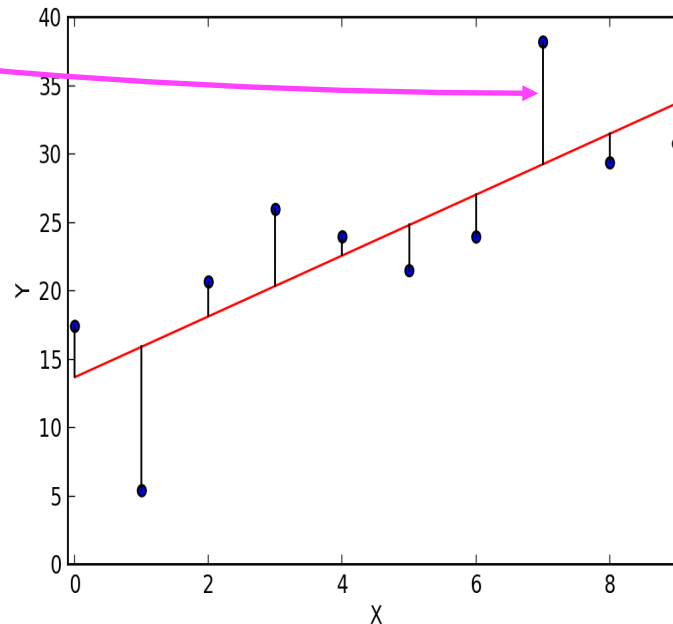# Method of least squares

● Error $= y - (\beta_0 + \beta_1 x_1)$

Sum of Squared Errors - $SSE = \sum n=1 \uparrow N \boxed{} (y \downarrow n - (B \downarrow 0 + B \downarrow 1 \; x \downarrow 1 + \dots + B \downarrow n \; x \downarrow n)) \uparrow 2$
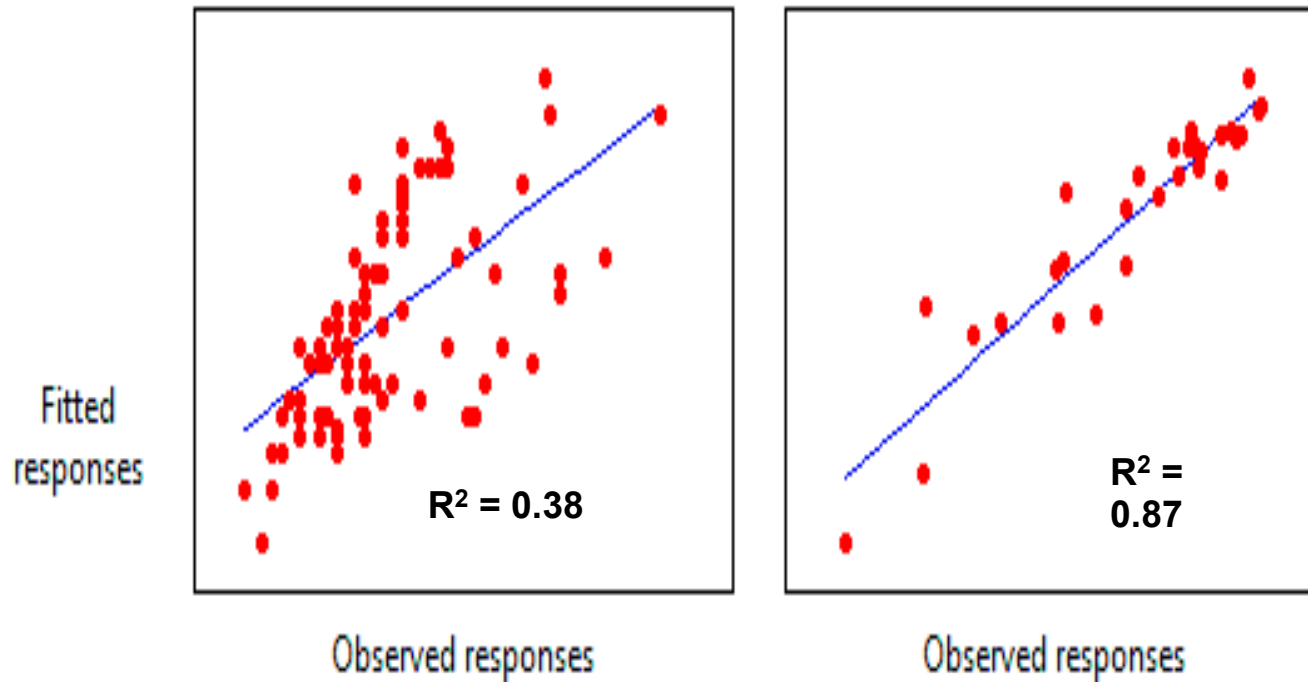
# R-squared

- How to measure the accuracy of the model?

- R-Squared – a statistical measure of how close the data are to the fitted regression line.

- R-squared = Explained variation / Total variation

- Between 0% - 100%

- 0 – none of the variation around the mean of the response variable data is explained

- 100 – all the variations accounted for

# R-Squared

**Plots of Observed Responses Versus Fitted Responses for Two Regression Models**

Fitted responses

$R^2 = 0.38$

$R^2 = 0.87$

Observed responses

Observed responses

http://blog.minitab.com/blog/adventures-in-statistics/regression-analysis-how-do-i-interpret-r-squared-and-assess-the-goodness-of-fit

# Recommender Systems

Wesley Schwerter, Shraddha Shridhar, Julio Butron, Juan Robledo

# How does a Recommender System Work

Two type of approaches
  1) Collaborative filtering
  ● Model based that uses either a single user's behavior or, more effectively, by using other users who have similar traits
  1) Content-based filtering
  ● Historical browsing information
  1) Hybrids

# Challenges with Recommender Systems

1) Some users can be modeled but other users do not exhibit typical behavior

- skewed results

- Decrease efficiency

2) Exploit a recommender system

- Favor one product over another

3) Scalability

- Performance

# The Netflix Recommender System

**NETFLIX**

- Streaming TV Shows and Movies is a fairly new concept
- Research shows that a typical Netflix subscriber checks on about 15 titles while browsing and may watch about 90 seconds of a title before they lose interest
- Goal is to present user with a number of videos they may like
- When DVDs were primary for Netflix, the star system was the focus for the recommender System
- Created a competition to improve accuracy of the star rating recommendation prediction.
- Use many algorithms for streaming
- Star ratings now shows you what they think you would rate the video

# The Netflix Recommender System

- Currently have billions of item ratings. Receive millions each day

- Millions of stream plays each day. Includes data like how long watched, time of the day, device type

- Metadata includes actor, director, ratings, reviews, etc….

- Records users amount of browsing, mouseovers, clicks, etc...

# Amazon Recommender System

- Before the dawn of video streaming s
  as
  Netflix and YouTube, Recommender
  were vital
  (and still are) to e-commerce.

- Recommender systems in e-commerce serve
  similar purposes: i.e. customize online store-front
  to each customer, bundle items that customer
  may need with initial purchase, and have physical
  inventory in appropriate regions in order to
  expedite shipping times.

# Challenges to E-commerce Recommender Systems

1) Amount of data can grow exponentially (Scalability)

2) Predicting needs for new customers is difficult due to limited information

3) E-commerce shoppers need to be able to find things quickly

4) Customer needs can be volatile and unpredictable

# Real world examples of Recommender Systems

Applications of Recommendation Systems

1. Product Recommendations
   a. Online retailers: Amazon.
2. Movies, Music Recommendations
   a. Netflix, Apple Music, Pandora, Spotify, Youtube, IMDB
3. News Recommendations
4. Online Dating Sites
5. Job/Career sites
6. Facebook

# What is big data?

The four V's

- Click stream
- Active/passive sensor
- Log
- Event
- Printed *corpus*
- Speech
- Social media
- Traditional

Volume

Variety

- Unstructured
- Semi-structured
- Structured

Big data

- Speed of generation
- Rate of analysis

Velocity

Veracity

- Untrusted
- Uncleansed

# Opportunities in the Education Domain

- Targeted interventions
- Faster and more in-depth diagnosis of learning needs
- Customizing learning based on each student's needs
- Institutions make better decisions
- Reduce dropout rates and Increases graduation numbers

# Conclusion

- **Big data**
  - Refers to analyzation of large amount of information collected through various channels
  - Significance in big data

- **How big data can help**
  - Designing instructional strategies
  - Feedback, motivation, personalization, efficiency, collaboration, tracking, understanding the learning process

- **How big data can hurt**
  - Privacy
  - Dehumanization
  - Deception by numbers

# Link Analysis Algorithms

- Page Rank
- Hubs and Authorities
- Topic-Specific Page Rank
- Spam Detection Algorithms
- Other interesting topics we won't cover
  - Detecting duplicates and mirrors
  - Mining for communities

# Ranking web pages

- Web pages are not equally "important"
    - www.joe-schmoe.com v www.stanford.edu

- Inlinks as votes
    - www.stanford.edu has 23,400 inlinks
    - www.joe-schmoe.com has 1 inlink

- Are all inlinks equal?
    - Recursive question!
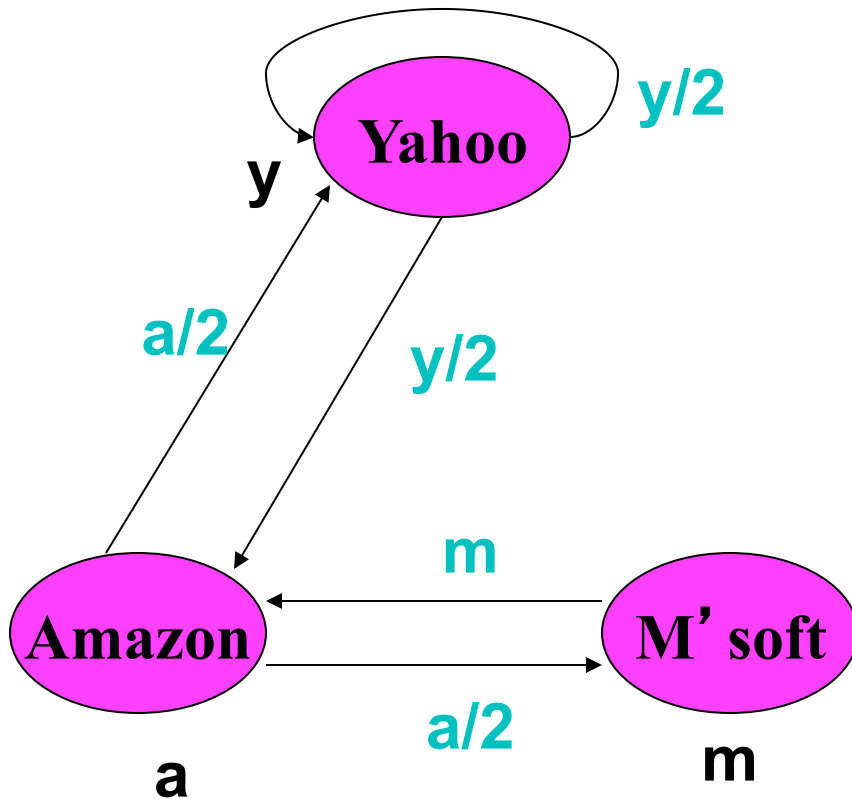
# Simple recursive formulation

- Each link's vote is proportional to the importance of its source page

- If page P with importance x has n outlinks, each link gets x/n votes

- Page P's own importance is the sum of the votes on its inlinks

# Simple "flow" model

The web in 1839



$y = y/2 + a/2$

$a = y/2 + m$

$m = a/2$